

Perceptual Models of Human-Robot Proxemics

Ross Mead and Maja J Matarić

Interaction Lab, Computer Science Department, University of Southern California,
3710 McClintock Avenue, RTH 423, Los Angeles, CA 90089-0781
{rossmead, mataric}@usc.edu
<http://robotics.usc.edu/interaction>

Abstract. To enable socially situated human-robot interaction, a robot must both understand and control *proxemics*—the social use of space—to employ communication mechanisms analogous to those used by humans. In this work, we considered how proxemic behavior is influenced by human speech and gesture production, and how this impacts robot speech and gesture recognition in face-to-face social interactions. We conducted a data collection to model these factors conditioned on distance. This resulting models of pose, speech, and gesture were consistent with related work in human-human interactions, but were inconsistent with related work in human-robot interactions—participants in our data collection positioned themselves much farther away than has been observed in related work. These models have been integrated into a situated autonomous proxemic robot controller, in which the robot selects interagent pose parameters to maximize its expectation to recognize natural human speech and body gestures during an interaction. This work contributes to the understanding of the underlying pre-cultural processes that govern human proxemic behavior, and has implications for the development of robust proxemic controllers for sociable and socially assistive robots situated in complex interactions (e.g., with multiple people or individuals with hearing/visual impairments) and environments (e.g., in which there is loud noise, reverberation, low lighting, or visual occlusion).

Keywords: Human-robot interaction, proxemics, Bayesian network

1 Motivation and Problem Statement

To facilitate face-to-face human-robot interaction (HRI), a *sociable* or *socially assistive robot* [1] often employs multimodal communication mechanisms similar 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). Like signals in electrical systems, these *social signals* are attenuated by distance; this influences how the signals are produced by people (e.g., humans adapt to increased distance by talking louder [2]), and subsequently impacts how these signals are perceived by the robot.

This research focuses on answering the following questions: Where do people position themselves when interacting with a robot? How does this positioning influence how people produce speech and gestures? How do positioning and human social signal production impact robot perception through automated speech and gesture recognition systems? How can the robot dynamically adjust its position to maximize its performance during the social interaction?

These questions are related to the field of *proxemics*, which is concerned with the interpretation, manipulation, and dynamics of spatial behavior in face-to-face social encounters [3]. Human-robot proxemics typically considers either a *physical* representation (e.g., distance and orientation [4–6]) or a *psychological* representation (e.g., amount of eye contact or friendliness [7, 8]) of the interaction. In [9], we proposed a probabilistic framework for *psychophysical* proxemic representation to bridge the gap between these physical and psychological representations by considering the perceptual experience of each agent (human and robot) in the interaction. We now formally investigate our framework proposed in [9], modeling position and perception to inform human-robot proxemics.

2 Related Work

The anthropologist Edward T. Hall [10] coined the term “proxemics”, and proposed that *cultural* norms define zones of intimate, personal, social, and public space [3]. These zones are characterized by the pre-cultural and *psychophysical* visual, auditory (voice loudness), olfactory, thermal, touch, and kinesthetic experiences of each interacting participant [2, 11] (Fig. 1). These *psychophysical* proxemic dimensions serve as an alternative to the sole analysis of distance and orientation (cf. [12]), and provide a functional perceptual explanation of the human use of space in social interactions. Hall [11] seeks not only to answer questions of *where* a person will be, but, also, the question of *why* they are there.

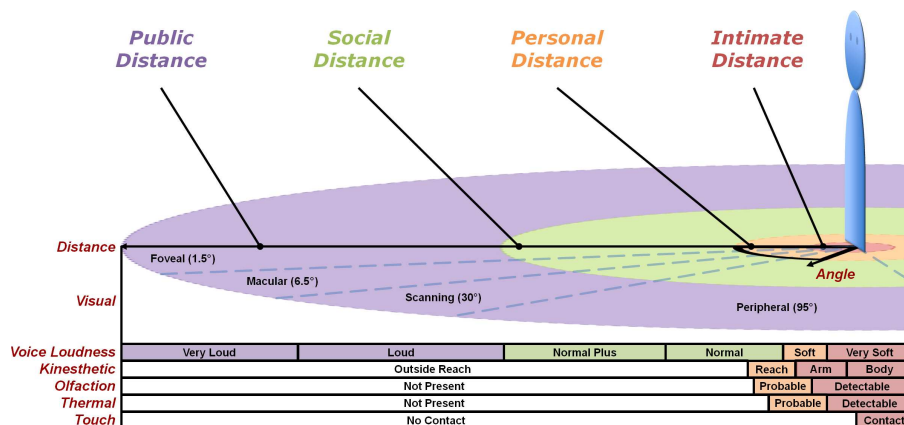


Fig. 1. Relationships between interpersonal pose and sensory experiences [3, 2, 11].

Contemporary probabilistic modeling techniques have been applied to socially appropriate person-aware robot navigation in dynamic crowded environments [13, 14], to calculate a robot approach trajectory to initiate interaction with a walking person [15], to recognize the averse and non-averse reactions of children with autism spectrum disorder to a socially assistive robot [16], and to position the robot for user comfort [4]. A lack of human-aware sensory systems has limited most efforts to coarse analyses [17, 18].

Our previous work utilized advancements in markerless motion capture (specifically, the Microsoft Kinect¹) 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 the physical [4–6] and the psychological [7, 8], considering psychophysical (perceptual) factors that contribute to both human-human and human-robot proxemic behavior [9, 23].

3 Framework for Human-Robot Proxemics

Our probabilistic proxemic framework considers how all represented sociable agents—both humans and robots—experience a co-present interaction [9]. We model the production (output) and perception (input) of speech and gesture conditioned on interagent pose (position and orientation).

3.1 Definition of Framework Parameters

Consider two sociable agents—in this case, a human (H) and a robot (R)—that are co-located and would like to interact.² At any point in time and from any location in the environment, the robot R must be capable of estimating:

1. an **interagent pose**, POS —Where will H stand relative to R ?
2. a **speech output level**, SOL_{HR} —How loudly will H speak to R ?
3. a **gesture output level**, GOL_{HR} —In what space will H gesture to R ?
4. a **speech input level**, SIL_{RH} —How well will R perceive H 's speech?
5. a **gesture input level**, GIL_{RH} —How well will R perceive H 's gestures?

These speech and gesture parameters are not concerned with the *meaning* of the behaviors, but, rather, the *manner* in which the behaviors are produced.

¹ <http://www.microsoft.com/en-us/kinectforwindows>

² Our current work considers only dyadic interactions (i.e., interactions between two agents); however, our framework is extensible [9], providing a principled approach (supported by related work [2]) for spatially situated interactions between any number of sociable agents by maximizing how each of them will produce and perceive speech and gestures, as well as any other social signals not considered by this work.

Interagent pose (POS) is expressed as a toe-to-toe distance (d) and two orientations—one from R to H (α), and one from H to R (β).

Speech output and input levels (SOL_{HR} and SIL_{RH}) are each represented as a sound pressure level, a logarithmic measure of sound pressure relative to a reference value, thus, serving as a signal-to-noise ratio. This relationship is particularly important when considering the impact of environmental auditory interference (detected as the same type of pressure signal) on SIL_{RH} and SOL_{HR} (our ongoing work).

Gesture output and input levels (GOL_{HR} and GIL_{RH}) are each represented as a 3D region of space called a gesture locus [24]. Related work in HRI suggests that nonverbal behaviors (e.g., gestures) should be parameterized based on proxemics [25]. For gesture production, we model the GOL_{HR} locus as the locations of H 's body parts in physical space. The GIL_{RH} can then be modeled as the region of R 's visual field occupied by the body parts associated with the gesture output of H (i.e., GOL_{HR}) [9].

3.2 Modeling Framework Parameters

We model distributions over these pose, speech, and gesture parameters and their relationships as a Bayesian network to represent: (1) how people position themselves relative to 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 (Fig. 2). Formally, each component of the model can be written respectively as:

$$p(POS) \quad (1)$$

$$p(SOL_{HR}, GOL_{HR} | POS) \quad (2)$$

$$p(SIL_{RH}, GIL_{RH} | SOL_{HR}, GOL_{HR}, POS) \quad (3)$$

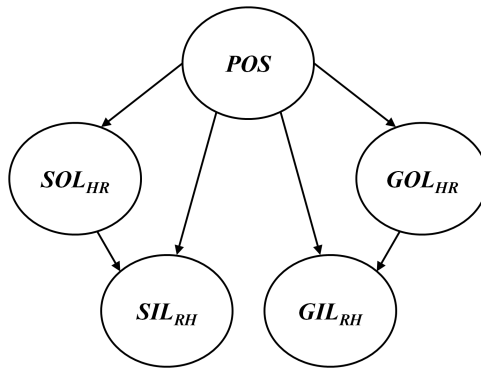


Fig. 2. A Bayesian network modeling relationships between pose, speech, and gesture.

4 Data Collection

We designed a data collection to inform the parameters of the models represented by Eqs. 1-3. We performed this experimental procedure in the context of face-to-face human-human and human-robot interactions at controlled distances.

4.1 Procedure

Each participant watched a short (1-2 minute) cartoon at separate stations (denoted $C1$ and $C2$ in Fig. 3), and then entered a shared space to discuss the cartoon with a social partner. Cartoons were chosen because they are commonly used in speech and gesture analysis studies, as people appear to gesture frequently when describing them [24]. Participants interacted for 4-6 minutes, and then separated into different rooms to watch another cartoon. Participants watched a total of six cartoons during the experimental session. For the first two (of six) interactions, participants were given no instructions as to where to stand in the room; this **“natural distance” within-participants condition** sought to elicit natural interagent poses, informing $p(POS)$ (i.e., Eq. 1), and also informing speech and gesture output/input parameters conditioned upon those natural poses (i.e., Eqs. 2 & 3). For the remaining four (of six) interactions, participants were instructed to stand at particular distances ($d = \{0.5, 2.0, 3.5, 5.0\}$ meters) with respect to their social partner (who stood at floor mark X in Fig. 3), the order of which was randomly selected for each interaction. This **“controlled distance” within-participants condition** sought to expose how people modulate their speech and gestures to compensate for changes in interagent pose, exploring the state space for the perceptual models represented in Eqs. 2 & 3.

For each participant, the social partner was either another participant or a PR2 robot³ (**“human vs. robot” between-participants condition**). For each interactions, participants watched and discussed three different cartoons and three of the same cartoons to explore the space of interaction dynamics. For human-robot interactions, participants were informed that the robot was learning how to communicate about cartoons. The robot generated speech and gestures based on annotated data selected from the human-human interactions.

4.2 Materials

Six Kinects mounted around the perimeter of the room were used to monitor interagent pose (POS), as well as to extract the positions of human body parts, which served as our representation of human gesture output levels (GOL_{HR}) [20]. The field-of-view of a Kinect on-board the robot allowed for the extraction of its gesture input levels (GIL_{RH}) based on tracked human body features (e.g., head, shoulders, arms, torso, hips, stance, etc.) [9]. Human speech output levels (SOL_{HR}) were recorded by calibrated microphone headsets; the corresponding robot speech input levels (SIL_{RH}) were recorded using the on-board Kinect microphone array.

³ <https://www.willowgarage.com/pages/pr2/overview>

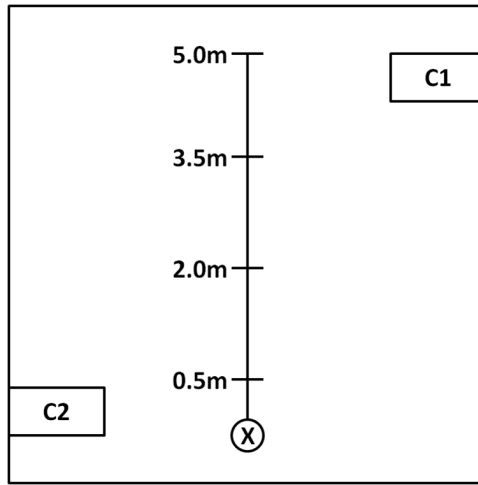


Fig. 3. The experimental setup for investigating proxemic behavior.

4.3 Participants

Participants were recruited via mailing list, flyers, and word-of-mouth on the campus of the University of Southern California; thus, these results might not apply in other regions (though the procedure would). A total of 40 participants (20 male, 20 female) were recruited for the data collection. All participants were university students between the ages of 18 and 35, and most had technical backgrounds; 10 of these participants had prior experiences with robots. None of the participants had ever interacted with each other prior to the experiment (i.e., they were “strangers” or, at most, “acquaintances”). Ethnicity was recorded (varied; predominantly North American), but not discriminated in the model, as we wanted to model framework parameters over a general population.

4.4 Dataset

We collected pose, speech, and gesture data for 10 human-human interactions (20 participants) and 20 human-robot interactions (20 participants). Each data collection session lasted for one hour, which included 20-45 minutes of interaction time. Recorded audio and visual data were independently annotated by two coders for occurrences of speech and gesture production; interrater reliability was high for both speech ($r_{speech} = 0.92$) and gesture ($r_{gesture} = 0.81$) annotations. For each of human-human and human-robot interactions, respectively, this dataset yielded 20 and 40 examples of natural distance selections, 4,914 and 4,464 continuous spoken phrases, and 2,284 and 1,804 continuous body gesture sequences; these numbers serve as the units of analysis in the following section.

5 Data Modeling and Analysis

The resulting dataset was used to model the relationships between speech and gesture output/input levels and interagent pose, as proposed in Eqs. 1-3.

5.1 Interagent Pose

To inform Eq. 1, interagent pose (*POS*) estimates from the Kinect during the **natural distance conditions** were used to generate the mean (μ) and standard deviation (σ) distance (in meters) in both human-human interactions ($\mu_{HH} = 1.44$, $\sigma_{HH} = 0.34$) and human-robot interactions ($\mu_{HR} = 0.94$, $\sigma_{HR} = 0.61$). An unpaired *t*-test revealed a very statistically significant difference ($p = 0.001$) in interagent pose under **human vs. robot conditions** (Fig. 4).

The interagent poses in our human-human interactions are consistent with human-human proxemics literature—participants in our data collection positioned themselves at a common social distance [3] as predicted by their interpersonal relationship⁴ [11].

However, the interagent poses in our human-robot interactions are inconsistent with human-robot proxemics literature—participants in our data collection positioned themselves much farther away than has been observed in related work. Walters *et al.* [6] consolidates and normalizes the results of many human-robot spacing studies, reporting mean distances of 0.49–0.71 meters under a variety of conditions. Takayama & Pantofaru [8] 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 two differences in experimental procedure. First, in many of these studies, participants are explicitly told to respond to a distance or comfort cue; however, in our study, participants were more focused on the interaction itself, so the positioning might have been less conscious. Second, in many of these studies, the robot is not producing gestures; however, in our study, the robot was gesturing and has a long reach (0.92 meters), so participants might have positioned themselves farther away from the robot to avoid physical contact. We will investigate this further in future work.

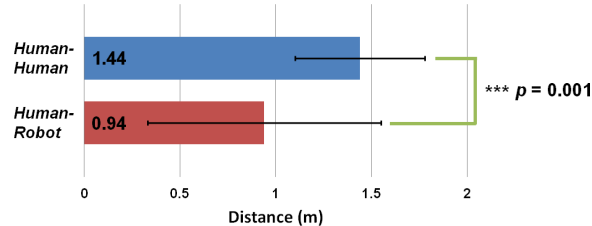


Fig. 4. Comparisons of human-human vs. human-robot proxemics ($p = 0.001$).

⁴ Recall that participants were “strangers” or, at most, “acquaintances”.

5.2 Human Speech/Gesture Output Levels

To inform Eq. 2, we considered human speech and gesture output levels (SOL_{HR} and GOL_{HR} , respectively). In the **natural distance condition**, human-human and human-robot speech output levels were ($\mu_{HH} = 65.80$, $\sigma_{HH} = 6.68$) and ($\mu_{HR} = 67.23$, $\sigma_{HR} = 5.90$) dB SPL, respectively. An unpaired t -test revealed no statistically significant difference between speech output levels under **natural distance, human vs. robot conditions**.

In the **controlled distance conditions**, an ANOVA F -test revealed significance ($p < 0.05$) in speech output levels at the nearest distance (0.5m) across the **human vs. robot conditions**. Furthermore, we found a trend towards significance ($p = 0.073$) for the statistical interaction of between distance and **human vs. robot conditions** (Fig. 5) for speech output levels; this suggests that the way in which a person modulates speech to compensate for distance might be different when interacting with a robot rather than a person. These results are consistent with [26], which found that people tend to speak more loudly with a robot when they think it is trying to recognize what they are saying, often compensating for a perceived lack of linguistic understanding by the robot; our work expands upon this by illustrating the effect at multiple distances.

We detected no significant difference in gesture output levels in any condition, so each body part is modeled as a uniform distribution within a person’s workspace (as measured by the Kinect). However, while our work suggests no significant relationship between gesture output level and *distance* (less than 5 meters) in dyadic face-to-face interactions, related work suggests that the *orientation* between people (e.g., L-shaped, side-by-side, and circular; [27]) influences the way in which people produce gestures [28]. For example, during a face-to-face interaction, we might expect to see gestures made directly in front of a person [28], thus, likely in the center of the field-of-view of the robot; however, during a non-frontal interaction, we expect to see gestures located more laterally [28], potentially falling out of the field-of-view of the robot. We will investigate the relationships between gesture output level and orientation in future work.

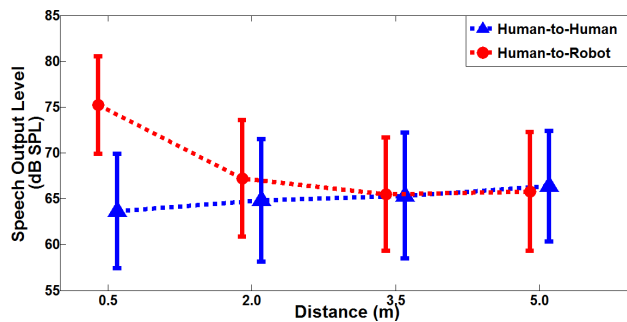


Fig. 5. Human speech output levels vary with distance in human vs. robot conditions.

5.3 Robot Speech/Gesture Input Levels

To inform Eq. 3, we used the data recorded by the microphones and camera of the Kinect on-board the robot to estimate continuous performance rates of automated speech and gesture recognition systems to inform our perceptual models of *SIL* and *GIL*, respectively (Fig. 6).

For speech recognition as a function of human-robot *POS*, we trained the model on annotated spoken phrases from a subset of the data (2 male, 2 female; 25 phrase each; 100 phrase vocabulary). Figure 6 illustrates the impact of Eq. 2 (estimating SOL_{HR} based on *POS*) on speech recognition rates (SIL_{RH}). The estimation of SOL_{HR} is important for accurately predicting system performance, especially because it often indicates that increase in vocal effort improves recognition; the alternative is to assume constant speaker behavior, which would predict a performance rate inversely proportional to the distance squared [29]).

For gesture recognition as a function of human-robot *POS*, we used the annotated gesture frames to calculate the number of times a particular body part appeared in the Kinect visual field versus how many times the body part was actually tracked (Fig. 7). Figure 6 illustrates the impact of Eq. 2 (estimating GOL_{HR} based on *POS*) on gesture recognition rates (GIL_{RH}) by estimating the joint probability of recognizing three body parts commonly used for human gesturing (in this case, the head and both hands) at different distances.

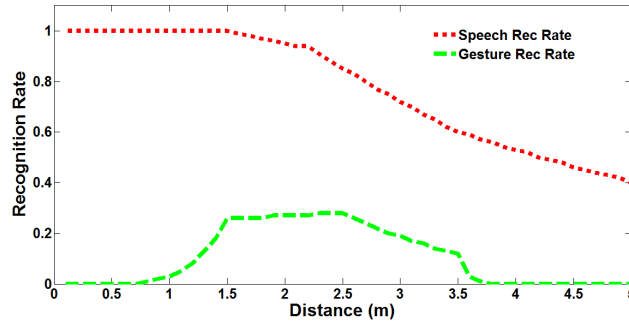


Fig. 6. Speech and gesture recognition rates as a function of distance.

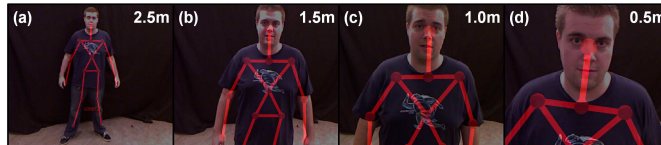


Fig. 7. Body features that fall into the Kinect field-of-view, depicted at four distances: (a) 2.5 meters, (b) 1.5m, (c) 1.0m, and (d) 0.5m. [9]

Our data collection revealed very little about the impact of interagent *orientation* (as opposed to *position*) on robot speech and gesture input levels, as face-to-face interaction does not promote exploration in the orientation space. Thus, we exercised alternative techniques to construct models relating orientation to SIL_{RH} and GIL_{RH} . As noted in Sec. 3.1, two orientations were considered: robot-to-human orientation (α) and human-to-robot orientation (β).

The impact of α on SIL_{RH} is modeled as the *head-related transfer function* (HRTF) of the PR2 robot listener (determined using the technique described in [30]). The impact of β is based on existing validated models of human *speaker directivity* [31]. The use of an HRTF and speaker directivity is commonly used in the audio processing community. Figure 8 illustrates the relationship between orientation and SIL_{RH} .

A small data collection was performed to model the relationship between orientation and gesture recognition rates. Participants stood at specified positions and orientations, and moved their limbs around in their kinematic workspace. GIL_{RH} was modeled based on human body features tracked by the Kinect [9] (Fig. 7). Figure 9 illustrates the relationships between orientation and GIL_{RH} .

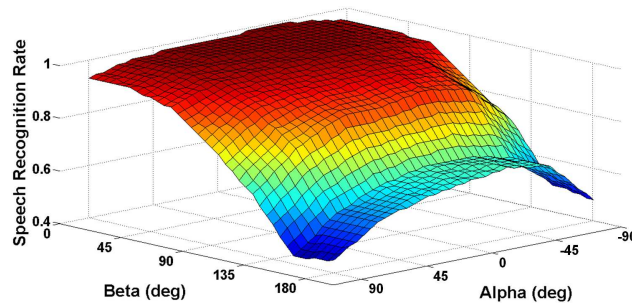


Fig. 8. Speech recognition rates as a function of speaker/listener orientations.

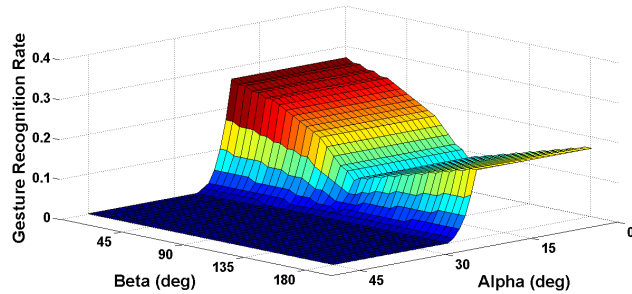


Fig. 9. Gesture recognition rates as a function of speaker/listener orientations.

6 From Perceptual Models to Autonomous Controllers

The fundamental insight of this work is in the application of these perceptual models of human-robot proxemics to enable spatially situated communication in HRI. Autonomous sociable robots must utilize automated recognition systems to reliably recognize natural human speech and body gestures. The reported models of human communication enable the system to predict the manner with which a behavior will be produced by a person (Sec. 5.2, Fig. 5), which can then be used to predict system performance (Sec. 5.3, Fig. 6). For example, if the robot can detect interagent pose, then it could use its models to predict (1) how loudly the user will likely speak, (2) how loudly its sensors will likely detect the speech, and (3) how well its speech recognition system is likely to perform—all before a single word is spoken by the person. If the sociable robot is mobile, it could use its predictions to inform a decision-making mechanism to decide to move to a better position to maximize the potential for its performance in the interaction. This is a fundamental capability that autonomous sociable robots should have, and has the potential to improve autonomy, increase richness of interactions, and generally make robots more reliable, adaptable, and usable.

In our ongoing work, we have integrated these models into a **situated autonomous proxemic robot controller**, in which the robot selects interagent pose parameters to maximize its expectation to recognize natural human speech and body gestures (Eq. 4). The controller utilizes a sampling-based approach, wherein each sample represents interagent distance and orientation, as well as estimates of human speech and gesture output levels (production) and subsequent robot speech and gesture input levels (recognition); a particle filter uses these estimates to maximize the expected performance of the robot during the interaction (Fig. 10).

As an extension of this, the robot can dynamically adjust its own speech and gesture output levels (SOL_{RH} and GOL_{RH} , respectively) to be consistent with the models of human-human interaction behavior reported in Sec. 5.2 (Eq. 5) to maximize human recognition of its own speech and gestures.

$$\operatorname{argmax}_{POS} E[SIL_{RH}, GIL_{RH} | SOL_{HR}, GOL_{HR}, POS] \quad (4)$$

$$\operatorname{argmax}_{SOL_{RH}, GOL_{RH}} E[SIL_{HR}, GIL_{HR} | SOL_{RH}, GOL_{RH}, POS] \quad (5)$$

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 pose, speech, and gesture parameters.

This controller has been implemented on the PR2 mobile robot (Fig. 11). The robot uses Eqs. 4 & 5 to estimate and select pose, speech, and gesture parameters to maximize its performance in the interaction, then uses the control equations described in our previous work [9] to realize these parameters.

We are in the process of validating this autonomous sociable robot controller.

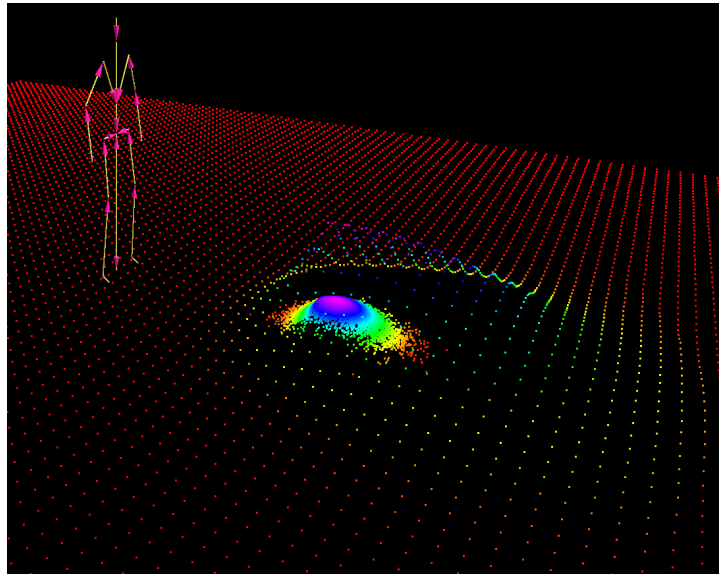


Fig. 10. Our autonomous sociable robot controller uses a sampling-based approach to estimate how the human produces social signals, and how the robot perceives them using models in our proxemic framework; a particle filter (with resampling) uses these estimates to maximize the expected performance of the robot during the interaction.

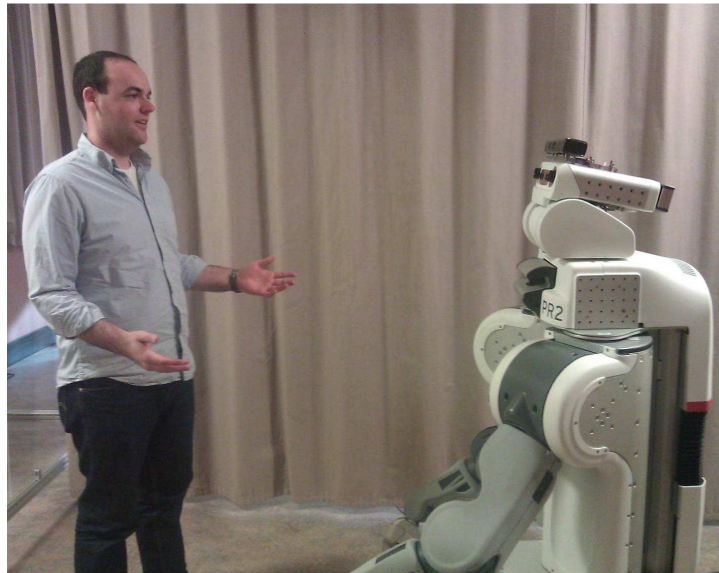


Fig. 11. The PR2 uses our controller to position itself in a conversation.

7 Summary and Contributions

In this work, we formally investigated our framework proposed in [9], modeling how proxemic behavior is influenced by human speech and gesture production, and how this impacts robot speech and gesture recognition in face-to-face social interactions. This resulting models of pose, speech, and gesture were consistent with related work in human-human interactions [2], but were inconsistent with related work in human-robot interactions [6, 8], warranting further investigation. This work contributes to the understanding of the underlying processes that govern human-human proxemic behavior [2], and provides an elegant extension into guiding principles of human-robot proxemic behavior. This research has implications for the development of robust spatially situated controllers for robots in complex interactions (e.g., with multiple people, or with individuals with hearing/visual impairments) and environments (e.g., in which there is loud noise, reverberation, low lighting, or visual occlusion).

The resulting models were used to implement a situated autonomous proxemic controller; the controller utilizes a sampling-based approach to maximize the expected performance of the robot during the interaction. We are in the process of validating the controller. 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 Repository⁵.

8 Future Work

Proxemic behavior is *not* cross-cultural [3], and the desired sensory experience of each interacting human participant varies all over the world [2]. 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 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 [1], 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 implementation of the system proposed in Sec. 6 has the robot producing speech and gesture based on models of how humans produce such behaviors (Sec. 5.2, Fig. 5). 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.

⁵ <https://code.google.com/p/usc-interaction-software>

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, dimly 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 model and integrate extrinsic interference into our probabilistic framework of human-robot proxemics.

Acknowledgments. This work is supported in part by an NSF Graduate Research Fellowship, the NSF National Robotics Initiative (IIS-1208500), NSF IIS-1117279 and CNS-0709296 grants, and the PR2 Beta Program.

References

1. Tapus, A., Matarić, M., Scassellati, B.: The grand challenges in socially assistive robotics. *IEEE Robotics and Automation Magazine* **14**(1) (2007) 35–42
2. Hall, E.T.: *The Hidden Dimension*. Doubleday Company, Chicago, Illinois (1966)
3. Hall, E.T.: *Handbook for Proxemic Research*. American Anthropology Association, Washington, D.C. (1974)
4. Torta, E., Cuijpers, R.H., Juola, J.F., van der Pol, D.: Design of robust robotic proxemic behaviour. In: *Proceedings of the Third International Conference on Social Robotics. ICSR'11* (2011) 21–30
5. Kuzuoka, H., Suzuki, Y., Yamashita, J., Yamazaki, K.: Reconfiguring spatial formation arrangement by robot body orientation. In: *HRI, Osaka, Japan* (2010)
6. Walters, M., Dautenhahn, K., Boekhorst, R., Koay, K., Syrdal, D., Nehaniv, C.: An empirical framework for human-robot proxemics. In: *New Frontiers in Human-Robot Interaction, Edinburgh* (2009) 144–149
7. Mumm, J., Mutlu, B.: Human-robot proxemics: Physical and psychological distancing in human-robot interaction. In: *HRI, Lausanne* (2011) 331–338
8. Takayama, L., Pantofaru, C.: Influences on proxemic behaviors in human-robot interaction. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems. IROS'09* (2009) 5495–5502
9. Mead, R., Matarić, M.J.: A probabilistic framework for autonomous proxemic control in situated and mobile human-robot interaction. In: *7th ACM/IEEE International Conference on Human-Robot Interaction, Boston, Massachusetts* (2012) 193–194
10. Hall, E.T.: *The Silent Language*. Doubleday Company, New York, New York (1959)
11. Hall, E.: A system for notation of proxemic behavior. *American Anthropologist* **65** (1963) 1003–1026
12. Mehrabian, A.: *Nonverbal Communication*. Aldine Transcation, Piscataway (1972)
13. Vasquez, D., Stein, P., Rios-Martinez, J., Escobedo, A., Spalanzani, A., Laugier, C.: Human aware navigation for assistive robotics. In: *Proceedings of the Thirteenth International Symposium on Experimental Robotics. ISER'12, Québec City, Canada* (2012)

14. Trautman, P., Krause, A.: Unfreezing the robot: Navigation in dense, interacting crowds. In: IROS, Taipei, Republic of China (2010)
15. Satake, S., Kanda, T., Glas, D.F., Imai, M., Ishiguro, H., Hagita, N.: How to approach humans?: Strategies for social robots to initiate interaction. In: HRI. (2009) 109–116
16. Feil-Seifer, D., Matarić, M.: Automated detection and classification of positive vs. negative robot interactions with children with autism using distance-based features. In: HRI'11, Lausanne, Switzerland (2011) 323–330
17. Oosterhout, T., Visser, A.: A visual method for robot proxemics measurements. In: HRI Workshop on Metrics for Human-Robot Interaction, Amsterdam (2008)
18. Jones, S., Aiello, J.: Proxemic behavior of black and white first-, third-, and fifth-grade children. *Journal of Personality and Social Psychology* **25**(1) (1973) 21–27
19. Mead, R., Atrash, A., Matarić, M.J.: Proxemic feature recognition for interactive robots: Automating metrics from the social sciences. In: International Conference on Social Robotics, Amsterdam, Netherlands (2011) 52–61
20. Mead, R., Atrash, A., Matarić, M.J.: Automated proxemic feature extraction and behavior recognition: Applications in human-robot interaction. *International Journal of Social Robotics* **5**(3) (2013) 367–378
21. Mead, R., Atrash, A., Matarić, M.J.: Recognition of spatial dynamics for predicting social interaction. In: 6th ACM/IEEE International Conference on Human-Robot Interaction, Lausanne, Switzerland (2011) 201–202
22. Feil-Seifer, D., Matarić, M.: People-aware navigation for goal-oriented behavior involving a human partner. In: IEEE International Conference on Development and Learning. Volume 2 of ICDL'11., Frankfurt Am Main, Germany (2011) 1–6
23. Mead, R., Atrash, A., Matarić, M.J.: Representations of proxemic behavior for human-machine interaction. In: NordiCHI 2012 Workshop on Proxemics in Human-Computer Interaction, Copenhagen, Denmark (2012)
24. McNeill, D.: *Hand and Mind: What Gestures Reveal about Thought*. Chicago University Press (1992)
25. Brooks, A.G., Arkin, R.C.: Behavioral overlays for non-verbal communication expression on a humanoid robot. *Autonomous Robots* **22**(1) (2007) 55–74
26. Kriz, S., Anderson, G., Trafton, J.G.: 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. HRI '10, IEEE Press (2010) 267–274
27. Kendon, A.: *Conducting Interaction - Patterns of Behavior in Focused Encounters*. Cambridge University Press, New York, New York (1990)
28. Ozyurek, A.: Do speakers design their co-speech gestures for their addressees? the effects of addressee location on representational gestures. *Journal of Memory and Language* **46**(4) (2002) 688–704
29. Wölfel, M., McDonough, J.: *Distant Speech Recognition*. John Wiley & Sons Ltd, West Sussex, United Kingdom (2009)
30. Gardner, B., Martin, K.: *Hrtf measurements of a kemar dummy-head microphone*. Technical Report 280, MIT Media Lab Perceptual Computing, Boston, Massachusetts (1994)
31. Chu, W., Warnock, A.: *Detailed Directivity of Sound Fields Around Human Talkers*. Research Report (Institute for Research in Construction (Canada)). Institute for Research in Construction (2002)