

# Space, Speech, and Gesture in Human-Robot Interaction

Ross Mead  
Interaction Lab  
University of Southern California  
Los Angeles, CA 90089-0781  
+1 213 740-6245  
rossmead@usc.edu

## ABSTRACT

To enable natural and productive 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: social speech and gesture production and recognition. My research focuses on answering these questions: How do social (auditory and visual) and environmental (noisy and occluding) stimuli influence spatially situated communication between humans and robots, and how should a robot dynamically adjust its communication mechanisms to maximize human perceptions of its social signals in the presence of extrinsic and intrinsic sensory interference?

## Categories and Subject Descriptors

I.2.9 [Artificial Intelligence]: Robotics – *operator interfaces*.

## General Terms

Algorithms; Design; Experimentation; Human Factors; Theory.

## Keywords

Human-robot interaction; proxemics; speech; gesture; multimodal.

## 1. INTRODUCTION

*If a person speaks or gestures, and no one is in a position to hear or see, is the person being social?*

*Proxemics* is the study of the interpretation, manipulation, and dynamics of spatial behavior in face-to-face social encounters. These dynamics are governed by sociocultural norms, which determine the overall sensory experience of social stimuli (speech, gesture, etc.) for each interacting participant [1]. Extrinsic and intrinsic sensory interference requires such spatially situated behavior to be dynamic. For example, if one is speaking in a quiet room, listeners need not 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.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICMI '12, October 22–26, 2012, Santa Monica, California, USA.  
Copyright 2012 ACM 978-1-4503-1467-1/12/10...\$15.00.

To facilitate natural and productive 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). My research focuses on answering the question: How do social (auditory and visual) and environmental (noisy and obstructive) stimuli 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 in the presence of extrinsic sensory interference (e.g., loud noise or visual occlusions) and intrinsic sensory interference (e.g., hearing or vision impairments)?

## 2. BACKGROUND

A number of rule-based spatial interaction controllers have been investigated in HRI [2-4]. Interpersonal dynamic models have been evaluated with sociable robots [5, 6]. Contemporary probabilistic modeling techniques have been applied to socially appropriate person-aware robot navigation in dynamic crowded environments [7], to calculating a robot approach trajectory to initiate interaction with a walking person [8], and to the recognition of averse and non-averse reactions of children with autism to a socially assistive robot [9].

My previous work utilized advancements in markerless motion capture (specifically, the Microsoft Kinect) to automatically extract low-level spatial features based on metrics from the social sciences [10]. These features were then used to recognize high-level spatiotemporal interaction behaviors, such as the initiation, acceptance, and termination of an interaction [11]. These investigations offered insights into the development of situated spatial interaction controllers for autonomous sociable robots, suggesting an alternative approach to proxemic behavior that goes beyond simple distance and orientation. Mead & Mataric [12] proposed a probabilistic framework for spatially situated interaction that considered the sensory experience of each agent (human or robot) in a co-present social encounter. This methodology established an elegant connection between previous approaches and illuminated the functional aspects of proxemic behavior in human-robot interaction, specifically, the impact of spacing on speech and gesture behavior recognition and production. My research is now aimed to formally model the relationship between these spatial, speech, and gesture parameters to inform spatially situated autonomous robot controllers for HRI.

### 3. APPROACH

My work will develop literature-grounded data-driven models of proxemic behavior that will be implemented as autonomous controllers for both telepresence and autonomous robots in HRI. Section 3.1 considers how all represented social agents—the person co-located with the robot (“co-located person”), the teleoperator (“remote person”), and the robot (telepresent or autonomous)—experience a co-present interaction; specifically, this work seeks to model the perception (input) and production (output) of auditory (speech) and visual (gesture) stimuli conditioned on interagent pose (position and orientation) and detected extrinsic interference (from the environment). Section 3.2 discusses a spatially situated autonomous controller for robot social behavior. A situated robot should be capable of generating a goal state (a set of spatial, speech, and gesture parameters) that is optimal for engaging in social interaction based on the estimated sensory experience of all interacting agents. The robot should then produce behaviors based on this goal state; this includes setting its speech and gesture output levels, as well as traversing a trajectory to an estimated goal pose, for which the model makes performance predictions about its speech and gesture recognition systems. Once the goal pose is reached, the robot should maintain the interaction, while responding to any changes in the interaction (as a result of movement, environmental interference, or additional co-located agents joining or leaving).

#### 3.1 Space, Speech, and Gesture Model

The parameters of the model are presented first, then the data that will be used to instantiate the model, noting that the data might not be available in a given interaction and thus should be inferred.

##### 3.1.1 Model Parameters

Social agents  $A$  and  $B$  are co-located in a room and would like to interact. Ideally, at any point in time and from any location in the room,  $A$  should be capable of estimating:

1. a **base pose**,  $POS_A$  — a position and orientation representing where it is located with respect to  $B$ ;
2. a **speech output level**,  $SOL_A$  — a sound pressure level representing how loudly it should speak for  $B$  to hear (from  $POS_A$ );
3. a **speech input level**,  $SIL_A$  — a sound pressure level representing how loudly it will likely perceive speech produced by  $B$  (from  $POS_A$ );
4. a **gesture output level**,  $GOL_A$  — a *gesture locus* [14] representing the space in which its gestures should occupy for  $B$  to see (from  $POS_A$ );
5. a **gesture input level**,  $GIL_A$  — a *gesture locus* [14] representing the space in which it will likely perceive gestures produced by  $B$  (from  $POS_A$ ).

In addition,  $A$  should be able to estimate how desirable these pose, speech input/output, and gesture input/output values are, based on sensed *extrinsic interference* ( $INT$ ) in the room (e.g., loud noise or visual occlusions) and should be able to estimate these values even if neither agent has spoken or gestured yet in the interaction.

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

##### 3.1.2 Data Collection

Mead et al. [10] describes the software and hardware infrastructure developed for data collection and visualization (Figure 1). The Microsoft Kinect will be used to monitor the relative pose between agents ( $POS$ ). Sound pressure levels ( $SIL$ ,  $SOL$ ) will be detected and produced by microphones and speakers, respectively. Mead et al. [10] demonstrated the ability to extract 6-DOF pose information from Kinect joint positions, affording the modeling of continuous distributions of  $GOL$  with respect to an agent frame (e.g., a distribution of hand locations with respect to the agent’s base pose); these gesture output levels can serve as parameters for an existing robot gesture generator, described in my previous work [14]. The Kinect field-of-view also allows 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 [12] (Figure 2).

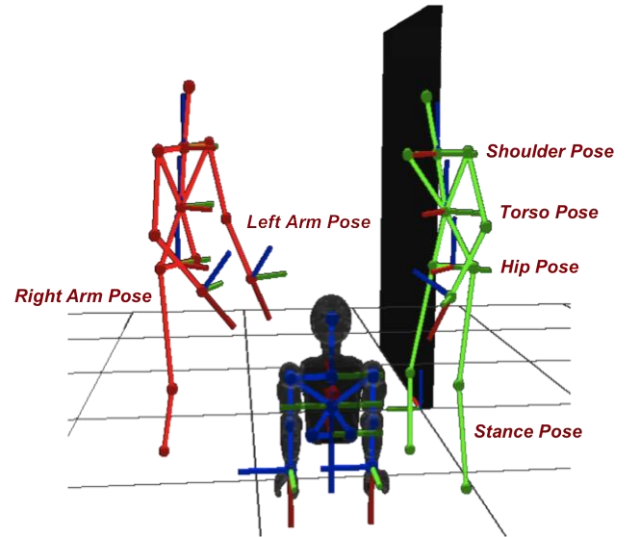


Figure 1: Noisy pose estimates of two humans and a robot interacting in a room with a visual occlusion (physical barrier) in the background. [10]

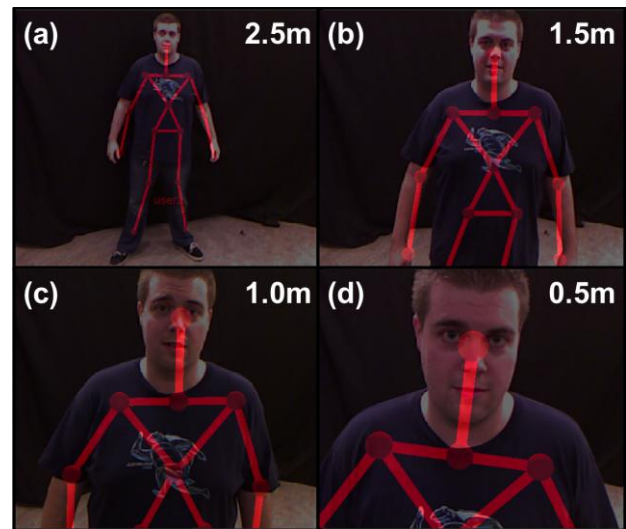


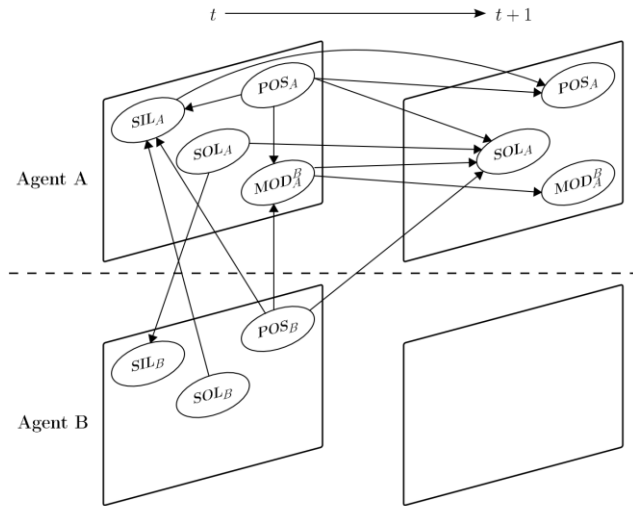
Figure 2: Body features that might fall into the visual field of a camera mounted on a robot, depicted at four distances: (a) 2.5 meters, (b) 1.5m, (c) 1.0m, and (d) 0.5m. [12]

Using the described sensor suites and agents (co-located person, remote person, and robot), a series of large-scale data collection sessions will be performed to expose how: (1) extrinsic interference influences interagent placement, (2) extrinsic interference and interagent placement influence agent speech and gesture production (output), and (3) extrinsic interference and interagent placement influence agent speech and gesture perception (input); formally:

1.  $p(POS_A | POS_B, INT)$
2.  $p(SOL_A, GOL_A | SIL_B, GIL_B, POS_A, POS_B, INT)$
3.  $p(SIL_A, GIL_A | SOL_B, GOL_B, POS_A, POS_B, INT)$

A co-located person will interact with either a telepresence or autonomous robot in simple interactions designed to highlight the discussed sensory factors. The telepresence condition will investigate how the remote person perceives and produces auditory and visual stimuli through a telepresence interface, and how these factors impact the manual positioning of the telepresence robot. The autonomous robot case will analyze how the robot's pose impacts the performance of its speech recognition and gesture recognition. Both scenarios will also collect auditory (from microphones) and visual data (from the Kinect) for estimating the experience of the co-located person. In all conditions, controlled amounts of extrinsic interference (loud noise, visual occlusions) will be introduced to determine their impact on interagent placement and social input/output levels.

It is acknowledged that a comprehensive data collection exercise for the large number of independent variables involved necessitates a large number of interaction scenarios. My work will focus on specific cases of interagent distance (close, medium, and far) and auditory and visual interference (no noise, small, medium, and large amounts of noise), and will motivate future work in adaptive models and controllers for situated HRI.



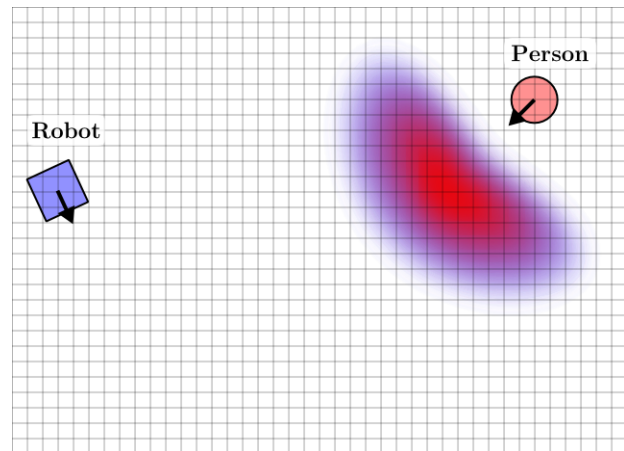
**Figure 3: The Dynamic Bayesian Network representation of the estimated perception and production of social stimuli; for clarity, only the base pose (POS) and speech input/output levels (SIL and SOL, respectively) are shown; the gesture input/output levels are modeled in the same way. Note that agent A also maintains a model of agent B ( $MOD_A^B$ ).**

### 3.1.3 Probabilistic Graphical Model

Parameters of distributions will be estimated over the collected data conditioned on the variables listed above; an agent A will maintain a model for itself (implicit), as well as a model for agent B ( $MOD_A^B$ ). These distributions and their interdependencies will be modeled as a Dynamic Bayesian Network [15] (Figure 3). At this point in the research, a co-located agent (telepresence or autonomous) would be capable of estimating its spatial, speech, and gesture parameters for engaging in a social interaction.

## 3.2 Space, Speech, and Gesture Controller

Using the Dynamic Bayesian Network (DBN), a spatially situated controller (for both telepresence and autonomous robots) could infer information about unobserved variables (e.g., the perceived speech and gesture input levels of a person) based on sensed data (e.g., interagent distance and environmental interference), and use these inferences to select parameters for robot social behavior (e.g., speech and gesture output levels) from some base pose. The robot will consider possible base poses that it could occupy, estimating how all interacting agents would perceive and produce social stimuli (speech and gesture) if the interaction occurred with the robot at that pose. As a first step, the space will be discretized using a grid-based approach (Figure 4). For each possible grid pose ( $POS_{x,y}$ ), the speech input/output levels ( $SIL_{x,y}$ ,  $SOL_{x,y}$ ) and gesture input/output levels ( $GIL_{x,y}$ ,  $GOL_{x,y}$ ) will be inferred from the DBN as if the interaction were to occur at that grid pose; the pose will then be assigned a weight ( $w_{x,y}$ ) based on the likelihood that an interaction would actually occur with these parameters under the sensed environmental conditions ( $INT$ ). The grid pose and output levels with the maximum weight will be selected as the goal state parameters; conditioning these parameters on the previous state ensures that they will not change drastically between timesteps (Figure 3). 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 an appropriate set of parameters.



**Figure 4: Pose and parameter likelihoods estimated by model.**

To traverse to the selected goal pose, the robot could then use a cost-based trajectory planner for safe navigation. Previous work by colleagues in the USC Interaction Lab has demonstrated the feasibility of real-time person-aware trajectory planning by weighting the calculated cost based on its fit to a more socially appropriate model of navigation [9]. The approach could be extended and applied to my research by weighting the trajectory cost by its fit to a model, such as that demonstrated in [11].

## 4. EVALUATION

To evaluate the developed models and resulting controllers, incremental as well as cumulative evaluation studies will be performed. First, individual models and the controller will be evaluated, then complete robot system implementations with the models in a variety of settings and tasks. To address generality, three different robot testbeds will be used, each capable of both telepresence and fully autonomous operation: the Willow Garage PR2, the Bandit upper-torso humanoid, and the Giraff mobile remote presence robot, all available in the USC Interaction Lab.

### 4.1 Individual Model/Controller Evaluations

The models of the auditory and visual sensory experiences of the co-located social partner, the remote social partner, and the robot, will each be individually evaluated using cross-validation on collected data. To evaluate the robot controller, two robot embodiments will be used (PR2 and Giraff; between participants), two modes of autonomy (manual and mixed initiative control; within-participants), and four cases of extrinsic interference (none, auditory, visual, and auditory+visual; within-participants). The evaluation studies will be performed at USC, with 40 participants across the adult age-span.

### 4.2 Complete System Evaluation

In a large multi-session study, 60 older adults from a partner senior living facility will each participate in four experimental sessions. All groups will interact with the Bandit robot platform in the context of seated exercises. In all cases, Bandit will enter the room, guide the participant(s) to the designated exercise chair, position itself in the coach location, and lead an exercise session. The first two of four sessions are one-on-one, the second two sessions are group sessions. To facilitate socialization among facility residents, this work aims to scale up to small groups of four participants, within the technical capabilities of the robot (perception, as well as interaction). Thus, the groups will consist of at least two or as many as four facility residents sitting in chairs arranged in a line and being coached by Bandit. Using group sessions will test the robot's spatially situated behavior, as well as any secondary socializing effects the session might have on the residents (i.e., do they discuss about the shared session, do they tend to socialize more as a result of it, etc.). This provides the opportunity to collect data and evaluate outcomes related to controller quality and effectiveness, one-on-one and in group settings, and overall task performance and adherence.

## 5. CONTRIBUTIONS

This research will contribute to the foundation of HRI, spanning both mobile remote presence and situated robot domains, providing validated, general, and reusable data-driven models and software for use by the robotics community at large. All code will be part of the Social Behavior Library (SBL) in the USC ROS Packages Repository (<http://sourceforge.net/projects/usc-ros-pkg>).

## 6. 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. I thank Maja Matarić (PhD advisor), Leila Takayama, Clifford Nass, Fei Sha, Louis-Philippe Morency, Amin Atrash, David Feil-Seifer, and Edward Kaszubski for their support on this research endeavor.

## 7. REFERENCES

- [1] Hall, E. T. 1974. *Handbook for Proxemic Research*. Society for the Anthropology of Visual Comm., Washington, D.C.
- [2] Shi, C., Shimada, M., Kanda, T., Ishiguro, H., and Hagita, N. 2011. Spatial formation model for initiating conversation. *Proc. of Robotics: Science and Systems*, Los Angeles, CA.
- [3] Kuzuoka, H., Suzuki, Y., Yamashita, J., and Yamazaki, K. 2010. Reconfiguring spatial formation arrangement by robot body orientation. *Proc. of the 5th ACM/IEEE Int'l. Conf. on Human-Robot Interaction*, Osaka, Japan, 285-292.
- [4] Walters, M. L., Dautenhahn, K., Te Boekhorst, R., Koay, K. L., Syrdal, D. S., and Nehaniv, C. L. 2009. An empirical framework for human-robot proxemics. *New Frontiers in Human-Robot Interaction*, 144-149.
- [5] Mumm, J. and Mutlu, B. 2011. Human-robot proxemics: physical and psychological distancing in human-robot interaction. *Proc. of the 6th Int'l. Conf. on Human-Robot Interaction*, Lausanne, Switzerland, 331-338.
- [6] Takayama, L. and Pantofaru, C. 2009. Influences on proxemic behaviors in human-robot interaction. *Proc. of the IEEE/RSJ Int'l. Conf. on Intelligent Robots and Systems*, St. Louis, MO, 5495-5502.
- [7] Trautman, P. and Krause, A. 2010. Unfreezing the robot: navigation in dense, interacting crowds. *Proc. of the IEEE/RSJ Int'l. Conf. on Intelligent Robots and Systems*, Taipei, Taiwan, 797-803.
- [8] Satake, S., Kanda, T., Glas, D. F., Imai, M., Ishiguro, H., and Hagita, N. 2009. How to approach humans?: strategies for social robots to initiate interaction. *Proc. of the 4th ACM/IEEE Int'l. Conf. on Human Robot Interaction*, San Diego, CA, 109-116.
- [9] Feil-Seifer, D. and Matarić, M. 2011. People-aware navigation for goal-oriented behavior involving a human partner. *Proc. of the IEEE Int'l. Conf. on Development and Learning*, 2, Frankfurt Am Main, 1-6.
- [10] Mead, R., Atrash, A., and Matarić, M. J. 2011. Proxemic feature recognition for interactive robots: automating metrics from the social sciences. *Proc. of the 2011 Int'l. Conf. on Social Robotics*, Amsterdam, Netherlands, 52-61.
- [11] Mead, R., Atrash, A., and Matarić, M.J. 2011. Recognition of spatial dynamics for predicting social interaction. *Proc. of the 6th ACM/IEEE Int'l. Conf. on Human-Robot Interaction*, Lausanne, Switzerland, 201-202.
- [12] Mead, R. and Matarić, M.J. 2012. A probabilistic framework for autonomous proxemic control in situated and mobile human-robot interaction. *Proc. of the 6th ACM/IEEE Int'l. Conf. on Human-Robot Interaction*, Boston, MA, 201-202.
- [13] Rossini, N. 2004. The analysis of gesture: establishing a set of parameters. *Gesture-based Comm. in Human-Computer Interaction*, LNCS, 2915, Springer, 463-464.
- [14] Mead, R., Wade, E.R., Johnson, P., St. Clair, A., Chen, S., and Matarić, M.J. 2010. An architecture for rehabilitation task practice in socially assistive human-robot interaction. *Proc. of The 19th IEEE Int'l. Symp. in Robot and Human Interactive Communication*, Viareggio, Italy, 404-409.
- [15] Koller, D. and Friedman, N. 2009. *Probabilistic Graphical Models*. MIT Press, Cambridge, MA.