

SITUATED PROXEMICS AND MULTIMODAL COMMUNICATION:
SPACE, SPEECH, AND GESTURE IN HUMAN-ROBOT INTERACTION

by

Ross Alan Mead

A Dissertation Presented to the
FACULTY OF THE GRADUATE SCHOOL
UNIVERSITY OF SOUTHERN CALIFORNIA
In Partial Fulfillment of the
Requirements for the Degree
DOCTOR OF PHILOSOPHY
(COMPUTER SCIENCE)

May 2016

Copyright 2016

Ross Alan Mead

Dedication

To anyone who has ever been wrong, you do not learn anything by being right.

Acknowledgments

The path to the Ph.D. is a long and treacherous road plagued by obstacles, doubt, and failure. I give my sincere thanks to those who aided me on this journey to success.

- First and foremost, to my Ph.D. advisor, Prof. Maja Mataric, thank you for your unending guidance and support. Your intellectual, technical, and strategic insights made me the researcher I am today. “Commit!”
- To my dissertation committee—Prof. Gaurav Sukhatme and Prof. Gigi Ragusa—and my qualifying examination committee—Prof. Fei Sha and Prof. Jonathan Gratch—thank you for elevating my work and guiding my research direction.
- To Prof. Jerry Weinberg, you inspired me to pursue robotics and graduate school, and are the reason I am where I am today. I am forever grateful to you.
- To Edward Kaszubski and Dr. Amin Atrash, you offered the greatest technical and intellectual considerations to my work. Thank you for your advice and friendship.
- To the members of the USC Interaction Lab, our shared ideas and conversations always helped me “get some gears” turning. You inspire excellence and innovation.

- To the 86 graduate, undergraduate, and high schools students whom I have had the honor of mentoring over the years, this work is a result of your efforts.
- To the staff of the KISS Institute for Practical Robotics (KIPR), my experiences in the Botball Educational Robotics Program provided the foundations of my interest in robotics. I am honored have been so welcomed into the KIPR family. Botball students continue to give me hope for the future.
- To my friends in Edwardsville and the STL, your continued interest and encouragement in my pursuits convinced me I was doing something good with my life.
- To my Jinriksha bandmates—Chadd Haselhorst, Tom Goodbrake, Andy Sogor, and Scott Bryant—you are the most creative people I know. Armada Forever.
- To my family, thanks for always being there and understanding when I was not.

Table of Contents

Dedication	ii
Acknowledgments	iii
List of Figures	viii
List of Tables	xii
Abstract	xiii
Chapter 1: Introduction	1
1.1 Overview	1
1.1.1 Problem Statement: Proxemics and Communication	1
1.1.2 Approach: Unified Models of Proxemics and Communication	3
1.2 Contributions	6
1.3 Outline	9
Chapter 2: Background and Related Work	10
2.1 Overview	10
2.2 Human-Robot Interaction (HRI)	11
2.2.1 Sociable Robotics	11
2.2.2 Socially Assistive Robotics (SAR)	12
2.3 Proxemics	12
2.3.1 Proxemics in Human-Human Interaction	13
2.3.2 Proxemics in Human-Robot Interaction	14
2.3.3 Representations for Proxemic Behavior Analysis	16
2.4 Summary	20
Chapter 3: Framework for Proxemics and Multimodal Communication	21
3.1 Overview	21
3.2 Heuristic Approach	23
3.2.1 Features of Individual Representation	23

3.2.2	Features of Physical Representation	26
3.2.3	Features of Psychological Representation	27
3.2.4	Features of Psychophysical Representation	28
3.2.5	Discussion	30
3.3	Data-Driven Approach	31
3.3.1	Definition of Framework Parameters	31
3.3.2	Modeling Framework Parameters	33
3.3.3	Discussion	34
3.4	Summary	34
Chapter 4: Modeling Proxemics and Multimodal Communication		35
4.1	Overview	35
4.2	Models of the Heuristic Approach	36
4.2.1	Modeling Features of Individual Representation	37
4.2.2	Modeling Features of Physical Representation	38
4.2.3	Modeling Features of Psychological Representation	38
4.2.4	Modeling Features of Psychophysical Representation	40
4.2.5	Discussion	44
4.3	Models of the Data-Driven Approach	45
4.3.1	Data Collection	45
4.3.2	Data Modeling and Analysis	48
4.3.3	Extension: Adaptation in Complex Environments	55
4.3.4	Discussion	57
4.4	Summary	58
Chapter 5: Implementation and Evaluation		59
5.1	Overview	59
5.2	Study 1: Human Behavior Recognition	60
5.2.1	Data Collection	60
5.2.2	Data Modeling, Analysis, and Results: Objective Measures	64
5.2.3	Discussion	68
5.3	Study 2: Robot Behavior Generation	69
5.3.1	Proxemic Behavior and Multimodal Communication Systems	69
5.3.2	System Analysis: Objective Measures	77
5.3.3	Discussion	84
5.4	Study 3: Human Acceptance of Robot Behaviors	85
5.4.1	Experimental Setup	86
5.4.2	Experimental Procedure and Measures	89
5.4.3	Experimental Conditions	94
5.4.4	Experimental Hypotheses	101
5.4.5	Participants	101
5.4.6	Data Analysis and Results: Behavioral Measures	102

5.4.7	Data Analysis and Results: Subjective Measures	109
5.4.8	Discussion	112
5.5	Summary	116
Chapter 6:	Summary and Conclusions	117
6.1	Contributions	117
6.2	Future Work: Adaptation in Complex Interactions	120
6.2.1	Adaptive Models of Human Speech Output Levels (SOL_{HR})	121
6.2.2	Adaptive Models of Human Speech Input Levels (SIL_{HR})	125
6.2.3	Discussion	126
6.3	Conclusions	126
	Bibliography	127

List of Figures

2.1	<i>Psychological</i> factors dictate the desired <i>psychophysical</i> (sensory) experience of each agent, which is manifested <i>physically</i> through the manipulation of space via change in position and orientation (Mead et al., 2012, 2013).	17
3.1	Individual pose features for two human users and an upper-body humanoid robot within a single framework; the absence of some features—such as the head, arms, or legs—signifies a pose estimate with low confidence.	25
3.2	In this interaction scenario, proxemic behavior is analyzed using simple physical features between each social dyad (pair of individuals).	26
3.3	Public, social, personal, and intimate distance codes, and SFP axis codes.	27
3.4	The anticipated sensory sensations that an individual would likely experience in different physical proxemic configurations and within the psychological distance zones.	28
3.5	Bayesian network modeling relationships between pose, speech, and gestures.	33
4.1	Error model of psychological distance code.	39
4.2	Error model of psychological sociofugal-sociopetal (SFP) axis code.	39
4.3	Error model of psychophysical visual code.	41
4.4	Error model of psychophysical voice loudness code.	41
4.5	Error model of psychophysical kinesthetic code.	42
4.6	Error model of psychophysical olfaction code.	42

4.7	Error model of psychophysical thermal code.	43
4.8	Error model of psychophysical touch code.	43
4.9	The experimental setup for modeling proxemic behavior and multimodal communication. Participants watched cartoons at locations $C1$ and $C2$. After watching each cartoon, one social partner (either human or robot; human vs. robot condition) relocated to the floor mark X . A participant approached the social partner along the line, and either a) stopped in at any interagent distance (natural distance conditions), or b) at one of four specified distances ($d = \{0.5, 2.0, 3.5, 5.0\}$ meters; controlled distance conditions).	46
4.10	A participant uses gestures while describing a cartoon to the PR2 robot.	46
4.11	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 (Mead and Matarić, 2012).	47
4.12	Comparisons of human-human vs. human-robot proxemics ($p = 0.001$).	50
4.13	Human speech output levels vary with distance in human vs. robot conditions.	51
4.14	Speech and gesture recognition rates as a function of distance.	53
4.15	Speech recognition rates as a function of human speaker orientation (β) and robot listener orientation (α).	54
4.16	Gesture recognition rates as a function of human speaker orientation (β) and robot observer orientation (α).	54
4.17	A Bayesian network modeling relationships between extrinsic (environmental) interference, pose, speech, and gesture.	56
5.1	The experimental setup for eliciting and recognizing human initiation and termination cues.	62
5.2	A five-state left-right HMM with two skip-states (Rabiner, 1990) used to model each modeled interaction cue— <i>initiation</i> or <i>termination</i> —for each representation.	65
5.3	Comparison of HMM classification accuracy of initiation and termination behaviors trained over physical, psychological, and psychophysical feature sets.	67

5.4	The goal state estimation system uses a sampling-based approach to estimate how the human produces social signals, and how the robot perceives them using models in the data-driven computational framework of proxemics and multimodal communication. Low to high estimates are denoted by both the coloring (red to magenta, respectively) and height (low to high, respectively) of sampled points in the floor plane (shown in black). The clustering of points denotes the impact of the resampling process. The system uses these estimates to select parameters that maximize the expected performance of the robot during the interaction.	72
5.5	The PR2 robot utilizes the reactive proxemic controller to reach the desired pose provided by the goal state estimation system during an interaction with a user.	74
5.6	Trajectories of an autonomous robot to a goal pose using three different motion control systems: 1) the reactive proxemic controller (red line; Section 5.3.1.2); 2) the <i>IP</i> -weighted cost-based trajectory planner (green line; Section 5.3.1.3); and 3) an unweighted cost-based trajectory planner (blue line) (Marder-Eppstein et al., 2010).	77
5.7	The experimental setup (left) and trajectories from the three motion control systems (right): 1) the reactive proxemic controller (red line), 2) the <i>IP</i> -weighted cost-based trajectory planner (green line), and 3) the unweighted cost-based trajectory planner (blue line).	80
5.8	A comparison of the <i>average normalized path length</i> (top) and <i>average interaction potential</i> (bottom) of each of the three motion control systems: 1) the reactive proxemic controller, 2) the <i>IP</i> -weighted cost-based trajectory planner, and 3) the unweighted cost-based trajectory planner.	82
5.9	The Bandit upper-body humanoid robot platform.	86
5.10	The experimental setup for evaluating human acceptance of robot behaviors.	87
5.11	The six phases of the experimental procedure.	89
5.12	Manipulation condition varying the maximum performance distance , $\mu_x = \{0.75, 1.25, 1.75, 2.25, 2.75, 3.25, 3.75, 4.25\}$ meters.	97
5.13	Manipulation condition varying the maximum performance value , $p_{max} = \{0.20, 0.40, 0.60, 0.80\}$	98
5.14	Manipulation condition varying the minimum performance value , $p_{min} = \{0.20, 0.40, 0.60, 0.80\}$	99

5.15	Baseline condition in which performance vs. distance is represented by a uniform distribution, $p(x) = p_{max} = p_{min} = 0.40$	100
5.16	Participant perceived location of robot peak performance (<i>perc</i>) vs. actual location of robot peak performance (<i>peak</i>). Note the heteroscedasticity of the data, which prevents us from performing traditional statistical analyses without first transforming the data (shown in Figure 5.17). . .	104
5.17	Participant perceived location of robot peak performance (<i>perc</i>) vs. actual location of robot peak performance (<i>peak</i>) on a log-log scale, reducing the effects of heteroscedasticity and allowing us to perform regression to determine parameters of the Power Law, ax^b	105
5.18	Changes in participant pre-/post-interaction proxemic preferences (<i>pre</i> and <i>post</i> , respectively; θ is the <i>contextual offset</i> defined in Section 5.4.6.1) vs. distance from participant pre-interaction proxemic preference (<i>pre</i>) to the actual location of robot peak performance (<i>peak</i>).	107
5.19	Changes in participant pre-/post-interaction proxemic preferences (<i>pre</i> and <i>post</i> , respectively; θ is the <i>contextual offset</i> defined in Section 5.4.6.1) vs. distance from participant pre-interaction proxemic preference (<i>pre</i>) to the perceived location of robot peak performance (<i>perc</i>).	108
5.20	The significant relationships modeled between the manipulated predictor variables and subjective measures. The correlation coefficient (β) and statistical significance (p) for each predictor-measure pair is presented along the connecting line; a dotted line indicates marginal significance. .	109
6.1	A graphical summary of the factors that influence the adaptation of human speech output levels (SOL_{HR}) over time. Red or blue indicates a state of either the robot or the human, respectively. A solid or dotted border indicates a state that is either measurable or latent (hidden), respectively.	124
6.2	A graphical summary of the factors that influence the adaptation of human speech input levels (SIL_{HR}) over time. Red or blue indicates a state of either the robot or the human, respectively. The grey rounded box indicates a “hearing impairment/sensitivity” (<i>HIS</i>) classification system. A solid or dotted border indicates a state that is either measurable or latent (hidden), respectively.	124

List of Tables

5.1	The observation vectors for 7-dimensional physical features, 3-dimensional psychological features, and 8-dimensional psychophysical features. . . .	65
5.2	Confusion matrix for HMM-based recognition of initiation and termination cues using physical, psychological, and psychophysical feature representations.	67
5.3	A comparison of the evaluation of distances predicted by two models— <i>robot performance</i> (using the goal state estimation system; Section 5.3.1.1) and <i>human preference</i> (reported in Section 4.3.2.1).	78
5.4	The distribution of affirmative responses provided by the robot across baseline (BL) and parameter-varying conditions ; note that average performance values p_{avg} vary as well. Manipulated values in each condition are highlighted in <i>bold italics</i>	96

Abstract

To facilitate face-to-face human-robot interaction (HRI), a *sociable robot* must employ 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 any other signals, these *social signals* are affected by distance and interference present in the medium through which they travel. People often compensate for this attenuation by adjusting the production of their social signals to compensate for these effects—for example, by speaking louder, using more broad gestures, or moving closer. How can a sociable robot do the same?

This dissertation investigates how *social* (speech and gesture) and *environmental* (loud noises and reduced visibility) factors influence positioning and communication between humans and sociable robots. Specifically, this research answers the following questions: 1) How should a robot dynamically adjust its position (proxemics) to maximize its automated recognition of human social signals? 2) How should a robot adjust its own communication behaviors to maximize human perceptions of its social signals? 3) How can a robot quickly adapt its models of proxemic and communication behavior to differences in human social signal perception?

This research formalizes *an extensible unifying framework for situated proxemics and multimodal communication in HRI*. The framework considers how both humans and robots experience social signals in face-to-face interactions. Data collections were conducted to inform probabilistic graphical models based on the framework that predict how speech and gesture are produced (transmitted) and perceived (received) by both humans and robots at different distances and under environmental interference.

This work integrates the resulting data-driven models into autonomous proxemic behavior and multimodal communication control system for sociable robots. The robot control system selects positioning parameters to maximize its ability to automatically recognize natural human speech and gestures. Furthermore, the robot control system can dynamically adjust its own speech and gestures to maximize human perceptions of its social signals. Experiments were conducted that successfully evaluated user acceptance of the autonomous robot proxemic control system, demonstrating that human users are willing to adapt their behavior preferences in exchange for improved robot performance in social contexts.

This research establishes a foundational component of HRI, enabling the development of robust controllers for *socially intelligent* robots in complex environments. Furthermore, this work has implications for technology personalization in *socially assistive* contexts with people with special needs, such as older adults, children with autism spectrum disorders, and people with hearing or visual impairments or sensitivities.

Chapter 1

Introduction

This chapter provides an overview of human-robot proxemics and multimodal communication. A novel representation and computational framework is introduced that relates proxemics and multimodal communication. Applications of the approach are discussed within the contexts of automated human behavior recognition and autonomous robot control systems, as well as implications for the design of robot social behaviors in face-to-face human-robot interactions. The chapter concludes with a list of primary and secondary contributions of this work, and an outline of the rest of the document.

1.1 Overview

1.1.1 Problem Statement: Proxemics and Communication

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 (Hall, 1959). These dynamics are governed by sociocultural norms, which determine the overall sensory experience of social stimuli (speech, gesture, etc.) for each interacting participant (Hall, 1974). Extrinsic (Adams and Zuckerman, 1991, Lloyd, 2009) and intrinsic (Hayduk and Mainprize, 1980, Malenby, 1975, Webb and Weber, 2003) sensory interference requires such spatial behavior to be dynamic. For example, if one is speaking in a quiet room, listeners need only be a few meters away to hear; however, if one is speaking in a loud 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 their locations strategically to perceive the speech and body language of the speaker in the same way.

To facilitate face-to-face human-robot interaction (HRI), a *sociable robot* (Breazeal, 2004) or *socially assistive robot* (Feil-Seifer and Matarić, 2005, Tapus et al., 2007) 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 any other signals, these *social signals* are *attenuated* (lose signal strength) based on distance, orientation, interference and other factors of the medium through which the signals travel. This attenuation influences how signals are *produced* by the *transmitting* sociable agent (human or robot; e.g., humans adapt to increased distance by talking louder (Hall, 1966, Traunmüller and Eriksson, 2000)), and subsequently impacts how these signals are *perceived* by the *receiving* sociable agent.

This dissertation investigates how *social* factors (speech and gesture) and *environmental* factors (loud noises and visual occlusions) influence proxemics and multimodal communication between co-present humans and sociable robots. Specifically, the research in this dissertation answers the following questions:

1. How should a robot dynamically adjust its proxemic behavior (position and orientation) to maximize its automated recognition of human social signals?
2. How should a robot adjust its own multimodal communication behaviors (speech loudness and gesture locations) to maximize human perceptions of its social signals?
3. How can a robot quickly adapt its models of proxemic and communication behavior in the presence of environmental interference?

1.1.2 Approach: Unified Models of Proxemics and Communication

This dissertation considers the role of and factors that contribute to proxemics and multimodal communication in social encounters between humans and robots. This work develops principled literature-grounded data-driven computational models for autonomous understanding (recognition) and use (control) of proxemic behavior to enable situated communication in HRI. Autonomous sociable robots are being used in data collections and evaluations of these models. Experiments conducted with these robots inform data-driven probabilistic models to estimate interaction parameters over the actual and desired sensory experience (e.g., loudness of voice or body features located in the visual field) of all agents involved (robot and human) in different proxemic scenarios (varying in distance and orientation) and in the presence of environmental interference (loud noises and visual occlusions). Each model is first validated individually, and then integrated for a systems-level validation to inform autonomous recognition and control

systems for robots. These autonomous behavior systems use the underlying predictive models, as well as online model adaptation, to enable a robot to dynamically determine appropriate parameters (e.g., body position, voice loudness, gesture locations) to conduct the interaction with a human user. This dissertation work provides experimentally validated robot behavior design principles, as well as parametric proxemic and multimodal communication models and software for situated HRI.

1.1.2.1 Behavior Representation and Computational Framework

This research draws upon insights from related work in human-human proxemic behavior analysis (Hall, 1963) to develop a novel representation and computational framework that unifies proxemic behavior and multimodal communication. The framework considers the *psychophysical* (sensory) experience of each agent (human and robot) in a co-present social encounter. Two approaches—*heuristic* and *data-driven*—are presented that go beyond traditional proxemic behavior representations, which focus on *physical* factors (distance and orientation) (Kuzuoka et al., 2010, Shi et al., 2011) or *psychological* factors (interagent relationship) (Mumm and Mutlu, 2011) that contribute to proxemic behavior. The methodology addresses the functional aspects of proxemic behavior in HRI, and provides a natural connection between previous approaches.

1.1.2.2 Feature Extraction and Human Behavior Recognition

To investigate different representations of proxemics, the dissertation develops software systems for the automated real-time extraction of salient features of proxemic behavior. The target features were represented in three ways: 1) a *physical* representation (distance and orientation), motivated by metrics commonly used in the social sciences for analyzing proxemic behavior (Mehrabian, 1972); 2) a *psychological* representation, characterized by the interpersonal relationship between two agents (Hall, 1966); and

3) a novel *psychophysical* representation, which considers the aural, visual, thermal, olfactory, and somatosensory experiences of each agent in a particular proxemic configuration (Hall, 1963, 1966). These proxemic feature extraction systems were evaluated during social interactions between two people and a humanoid robot in the presence of a visual occlusion. The three proxemic feature representations are compared by training probabilistic graphical models (Koller and Friedman, 2009) to recognize human spatiotemporal behaviors that signify transitions into (initiation) and out of (termination) a social interaction. The models trained on the psychophysical features (encoding the sensory experience of each interacting agent) were shown to outperform those trained on traditional physical and psychological features, suggesting a more powerful representation of proxemic behavior in both human-human and human-robot interactions.

1.1.2.3 Predictive Models and Robot Behavior Generation

Insights gained from the development of a representation and feature extraction technique for proxemic behavior led to the consideration of how human social signal production (e.g., speech and gesture) is influenced by proxemic behavior, and how this impacts autonomous robot social signal recognition in face-to-face social interactions. This dissertation models how human users adapted their multimodal communication behaviors conditioned on interagent distance and environmental interference in both human-human and human-robot interactions. The resulting models were integrated into a situated autonomous proxemic robot controller, in which the robot selects interagent pose parameters to maximize its expectation (based on predictive models) to recognize natural human speech and body gestures during a face-to-face social interaction. Furthermore, the robot dynamically adjusts parameters of its own speech and gesture production (voice loudness and gesture locations, respectively) to be consistent with the models of human-human interaction behavior in an effort to maximize human

recognition of its own social signals. Environmental interference (e.g., loud noises or visual occlusions) influences how both the robot and the human produce and perceive social signals, necessitating situated adaptive models characterized by unique challenges and solutions.

1.1.2.4 System Evaluations with Human Users

The performance of the automated human behavior recognition and autonomous robot behavior generation systems developed in this work are evaluated with respect to objective, subjective, and behavioral measures in some of the largest studies ever conducted in the field of human-robot interaction. The analysis of this work a) demonstrates accurate predictions of human social behavior and robot recognition of these behaviors (objective evaluation), b) provides strong relationships between human perceptions of the robot based on its ability to recognize human social signals (subjective evaluation), and c) demonstrates that human users will naturally adapt their own interaction behaviors to improve robot performance (behavioral evaluation). These results have significant implications for the design and deployment of sociable robots in the real world.

1.2 Contributions

This dissertation contributes to the understanding of the underlying processes that govern human-human proxemic behavior, and provides a natural extension into guiding principles of human-robot proxemic behavior, establishing a foundational component of HRI. This research has implications for the development of robust and adaptive social behavior control systems for sociable robots situated in complex environments (e.g., in which there are loud noises or visual occlusions) and interactions (e.g., with multiple people, or with individuals with hearing or visual impairments).

The following are the primary contributions of this work:

1. *An extensible unifying framework for situated proxemics and multimodal communication for both human-human and human-robot interactions.* The framework considers how both humans and robots experience social signals in face-to-face interactions. Data collections were conducted to inform probabilistic graphical models that predict how speech and gesture are produced (transmitted) and perceived (received) by both humans and robots at different distances and under environmental interference.
2. *Proxemic feature extraction and behavior recognition systems.* The system automatically extracts proxemic features based on three feature representations—physical, psychological, and psychophysical. These features are used to recognize transitions into (initiation) and out of (termination) co-present social interactions. A comparison of representations is provided, in which the psychophysical representation of proxemic behavior outperforms traditional physical and psychological representations, suggesting a more powerful approach to the recognition of spatiotemporal interaction cues in both human-human and human-robot interactions.
3. *Probabilistic models for autonomous robot generation of proxemic behavior and multimodal communication.* Using the computational framework in this dissertation, proxemic and communication behavior are unified in a probabilistic graphical model that represents the production and perception of social signals (e.g., speech and gesture) as a function of interagent pose (distance and orientation). Pose, speech, and gesture parameters are selected to maximize social signal recognition rates for all agents (human and robot) in the interaction.

4. *A method for adapting robot proxemic and multimodal communication parameters in complex environments.* The approach has implications for situated robot behavior generation in complex environments (in which there are loud noises or visual occlusions), and for technology personalization in complex interactions with a focus on socially assistive contexts with people with special needs, such as those with hearing or visual impairments or sensitivities.
5. *Objective, subjective, and behavioral evaluations of proxemic and multimodal communication control systems.* Experiments were conducted that a) demonstrated accurate predications of human social signal production and robot social signal recognition (objective measures), b) related user perceptions of the robot to its ability to recognize social signals (subjective measures), and c) demonstrated that human users naturally adapt their own proxemic preferences to improve robot performance in social contexts (behavioral measures).

This dissertation also provides software and a corpus of public human-human and human-robot interaction data for use by researchers worldwide, serving to inform, validate, and extend longstanding research in the both HRI and the social sciences.

The following are the secondary contributions of this work:

1. *Implemented and validated open-source software systems.* The automated human behavior recognition and autonomous robot behavior control systems are publicly available in the Social Behavior Library (SBL), an open-source software suite that provides generic computational models of social behavior for HRI.
2. *Large corpora of human-human and human-robot interaction data.* This dissertation provides uniquely comprehensive corpora of human-human and human-robot interaction data, which are publicly available upon request to facilitate research in both robotics and social science communities.

1.3 Outline

The remainder of this document is organized as follows:

- Chapter 2 provides background on existing work in both human-human and human-robot proxemics, and the relationship between proxemics and multimodal communication in face-to-face social encounters.
- Chapter 3 describes two approaches—heuristic and data-driven—employed for unifying proxemics and multimodal communication into a computational framework for human-robot interaction.
- Chapter 4 details the process of modeling proxemics and multimodal communication using the heuristic and data-driven approaches.
- Chapter 5 discusses the evaluation of the models of proxemics and multimodal communication in the contexts of automated human behavior recognition (using the heuristic approach) and autonomous robot behavior generation (using the data-driven approach). Objective, subjective, and behavioral evaluations are also provided.
- Chapter 6 summarizes the contributions of the dissertation, and discusses open challenges and future work for extending the models.

Chapter 2

Background and Related Work

This chapter reviews the field of human-robot interaction, with an emphasis on autonomous sociable robots and socially assistive robots. Proxemic behavior is surveyed in both human-human and human-robot interactions. Common representations of proxemic behavior are identified, and a novel representation based on sensory experience is presented that relates proxemics and multimodal communication, providing a foundation for the approach and contributions of this dissertation.

2.1 Overview

This dissertation develops a computational framework unifying proxemic behavior and multimodal communication in human-robot interaction. Support for this work comes from related literature in autonomous HRI—specifically, sociable robotics and socially assistive robotics (reviewed in Sections 2.2.1 and 2.2.2, respectively)—as well as from literature on proxemics in both human-human and human-robot interactions (reviewed in Sections 2.3.1 and 2.3.2, respectively). Insights gained from related work suggested three categories of proxemic behavior representation (presented in Section 2.3.3), which provide the inspiration and foundation for the approaches of this dissertation.

2.2 Human-Robot Interaction (HRI)

The field of human-robot interaction (HRI) constitutes a broad range of interactions and application domains with human users, including human control of robots—such as in mobile remote presence (Lee and Takayama, 2011)—to autonomous robots that interact with people using natural communication, such as speech and gestures (Breazeal, 2004). This dissertation focuses on the latter. The sections below survey the fields of sociable robotics and socially assistive robotics.

2.2.1 Sociable Robotics

Dautenhahn (2007) proposes three categorical approaches to the investigation of social interactions with robots: (1) *human-centered HRI*, in which robot activity is evaluated based on its acceptance and comfort as perceived by a human; (2) *robot-centered HRI*, in which a robot has its own drives or motivations, some of which can be realized through interactions with a human; and (3) *robot cognition-centered HRI*, in which the processes of robot learning and decision-making can be influenced by interactions with a human. These categories are not mutually exclusive. Much of the work on human-robot interaction and human-robot proxemics has explored human-centered HRI (Section 2.3.2); however, Breazeal (2003) notes that “it is important to view the design and evaluation problem from the robot’s perspective as well as that of the human.” Such an approach falls between the categories of robot-centered and human-centered HRI, which is referred to as *socially situated HRI* (Dautenhahn, 2007). This is the approach taken by this dissertation in the development of computational framework for proxemic behavior and multimodal communication in HRI (Chapter 3).

2.2.2 Socially Assistive Robotics (SAR)

Socially assistive robotics (SAR) is the intersection of social robotics and assistive robotics that focuses on non-contact human-robot interaction aimed at monitoring, coaching, teaching, training, and rehabilitation domains (Feil-Seifer and Matarić, 2005). Notable areas of SAR include robotics for older adults (Fasola and Matarić, 2010, Libin and Cohen-Mansfield, 2004, Tapus et al., 2009a,b, Wada and Shibata, 2007, Wada et al., 2004), for children with autism spectrum disorders (Dautenhahn et al., 2009, Feil-Seifer and Matarić, 2008, 2009, 2011a, Kozima et al., 2007, Ricks and Colton, 2010, Scassellati, 2007, Welch et al., 2010), and for people in post-stroke rehabilitation (Matarić et al., 2007, Mead et al., 2010), among others. SAR systems have been shown to increase user motivation and engagement in exercises (Kidd and Breazeal, 2011, Wada et al., 2004), improve adherence to prescribed lifestyle (Fasola and Matarić, 2010, 2011), lower stress (Wada et al., 2004), and stimulate verbalization (Tapus et al., 2009a, Wada and Shibata, 2007) and socialization (Kim et al., 2008, Wainer et al., 2010). Consequently, SAR constitutes an important subfield of robotics with significant potential for health and quality of life. Because the majority of SAR contexts experimented with to date involve one-on-one interaction between the robot and the user, the goals of this dissertation are especially relevant, as they provide principled models proxemics and multimodal communication for such interactions, in the SAR contexts and beyond.

2.3 Proxemics

There exists a rich body of work in the social sciences that seeks to explain proxemic phenomena in human-human interactions (discussed in Section 2.3.1). Many of the factors that influence human-human proxemics have also been investigated in human-robot interactions (discussed in Section 2.3.2).

In Section 2.3.3, previous approaches for proxemic behavior analysis are grouped into three novel categorizations of representation for proxemic behavior: *physical*, *psychological*, and *psychophysical*. The psychophysical representation serves to unify proxemics and multimodal communication, discussed in Section 2.3.3.3.

2.3.1 Proxemics in Human-Human Interaction

The anthropologist Edward T. Hall (1959) coined the term proxemics, defining it as “the interrelated observations and theories of man’s use of space as a specialized elaboration of culture” and proposing four culture-specific zones of proxemic distance: public, social, personal, and intimate (Hall, 1966). Mehrabian (1972), Argyle and Dean (1965), and Burgoon et al. (1995) analyzed proxemic behavior by considering psychological indicators of the interpersonal relationship between social partners (e.g., amount of mirroring, reciprocity, and compensatory behaviors; amount of eye gaze and smiling; posture and arm configurations; intimacy of topic and thought; etc.).

Schöne (1984) was inspired by the spatial behaviors of biological organisms in response to stimuli, and investigated human spatial dynamics from the physiological and ethological perspectives of the human sensory system; similarly, Hayduk and Mainprize (1980) and Mallenby (1975) analyzed the personal space requirements of people with visual and hearing impairments, respectively. Kennedy et al. (2009) studied the amygdala and how fight-or-flight responses are involved in the regulation of interpersonal space. Kendon (1990) analyzed the organizational patterns of social encounters, categorizing them into so-called *F-formations*: “when two or more people sustain a spatial and orientation relationship in which the space between them is one to which they have equal, direct, and exclusive access”.

Schegloff (1998) proposed that people use body orientation, such as stance, hip and shoulder orientation, head pose, and eye gaze, to communicate an interest in initiating, accepting, maintaining, terminating, or avoiding an interaction (Deutsch, 1977). McNeill (2005) studied how people manipulate space in an interaction, either to direct attention to an external stimulus (e.g., using a pointing gesture with the hand) or to guide a social partner to a location (e.g., using spatial configurations similar to the F-formations reported by Kendon (1990)).

Human proxemic behavior is also impacted by factors of the individual—such as gender (Price and Dabbs Jr., 1974), age (Aiello, 1987), ethnicity (Jones and Aiello, 1973), and personality (Aiello, 1987)—as well as factors of the environment—such as lighting (Adams and Zuckerman, 1991), setting (Geden and Begeman, 1981), location in setting and crowding (Evans and Wener, 2007), size (Aiello et al., 1981), and permanence (Hall, 1966).

2.3.2 Proxemics in Human-Robot Interaction

The emergence of virtually and physically embodied conversational agents (including sociable robots) (Cassell et al., 2000) necessitated formal computational models of social proxemics in human-agent interactions. Many rule-based proxemic behavior controllers have been implemented for HRI (Hüttenrauch et al., 2006, Kuzuoka et al., 2010, Shi et al., 2011, Walters et al., 2009). Kirby et al. (2007, 2009) demonstrated socially-acceptable navigation strategies for person following and hallway pass-by behavior. Interpersonal dynamic models, such as equilibrium theory (Argyle and Dean, 1965), have been implemented and evaluated in HRI (Mumm and Mutlu, 2011, Takayama and Pantofaru, 2009).

Contemporary probabilistic modeling techniques have been applied to socially-appropriate person-aware robot navigation in dynamic crowded environments (using Gaussian processes) (Trautman and Krause, 2010), to calculating a robot approach trajectory to initiate interaction with a walking person (using support vector machines) (Satake et al., 2009), to the recognition of averse and non-averse reactions of children with autism spectrum disorders with respect to a socially assistive robot (using Gaussian mixture models) (Feil-Seifer and Matarić, 2011a), and to position the robot for user comfort (using a sampling-based approach) (Torta et al., 2011). A lack of high-resolution metrics and sensor systems limited previous efforts to coarse analyses in both space and time (Jones and Aiello, 1973, Oosterhout and Visser, 2008). Recent developments in markerless motion capture, such as the Microsoft Kinect¹, have addressed the problem of real-time human pose estimation, providing the means and justification to revisit and more accurately model the subtle dynamics of proxemic interaction.

The majority of proxemics work in HRI focuses solely on meeting the needs and preferences of a human user during a face-to-face interaction. The results of many human-robot proxemics studies are consolidated and normalized in Walters et al. (2009), reporting mean distances of 0.49–0.71 meters using a variety of robots and conditions. Preferences between humans and the PR2 robot² were investigated by Takayama and Pantofaru (2009), reporting mean distances of 0.25–0.52 meters; this dissertation investigated human proxemic preferences using the same PR2 robot platform, but in a different context (Section 4.3.2.1), reporting a mean distance preference of 0.94 meters (Mead and Matarić, 2014), illustrating the variability of these preferences. Farther proxemic preferences have been measured in Mumm and Mutlu (2011) and Torta et al. (2013), reporting mean distances of 1.0–1.1 meters and 1.7–1.8 meters, respectively.

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

²<https://www.willowgarage.com/pages/pr2/overview>

However, initial investigations of human-robot proxemics performed in this dissertation suggested immediate technical challenges for autonomous sociable robots using human preference-based proxemic configurations: at these distances, the robot did not perform well, as measured by automated speech and gesture recognition rates (Mead and Matarić, 2014). Speech recognition performed adequately at distances less than 2.5 meters, and face and hand gesture recognition performed well at distances of 1.4–2.5 meters; thus, given current technologies, distances for mutual recognition of these social signals is between 1.4 and 2.5 meters, at and beyond the far bounds of previously reported human proxemic preferences (presented formally in Section 4.3.2.3). This insight led to the formalization of representations used in proxemic behavior analysis, in an effort to unify human-robot proxemics and perception of multimodal communication, discussed in the next section.

2.3.3 Representations for Proxemic Behavior Analysis

This dissertation consolidates previous work in proxemic behavior analysis into three related categories of representation for proxemic behavior:

1. the **physical representation**, based on distance and orientation (introduced in Section 2.3.3.1, and formalized in Section 3.2.2);
2. the **psychological representation**, based on the interpersonal relationship (introduced in Section 2.3.3.2, and formalized in Section 3.2.3); and
3. the **psychophysical representation**, based on the sensory experience of social stimuli (introduced in Section 2.3.3.3, and formalized in Section 3.2.4).

2.3.3.1 Physical Representation of Proxemic Behavior

The *physical* representation of proxemic behavior is concerned with how space is occupied by two or more bodies (Hall, 1959, Kuzuoka et al., 2010), relating these bodies via low-level spatial parameters of distance and orientation (Mehrabian, 1972, Schegloff, 1998) (Figure 2.1). It is the most commonly used representation in the analysis and manipulation of proxemics in both human-human and human-machine interactions. Kastanis and Slater (2012) manipulated physical parameters of a virtual human to predictably influence the position of real people in immersive virtual environments. Marquardt and Greenberg (2012) utilized physical proxemic features over time (i.e., movement), and relative to specific people or objects in an environment, for context-aware ubiquitous computing scenarios. Many control architectures in mobile HRI rely on these physical parameters alone for autonomously controlling human-robot proxemic configurations (Hüttenrauch et al., 2006, Kuzuoka et al., 2010, Walters et al., 2009).

2.3.3.2 Psychological Representation of Proxemic Behavior

The *psychological* representation of proxemic behavior is concerned with the high-level interpersonal relationship between two or more agents (Hall, 1966) (Figure 2.1). For example, Affiliative Conflict Theory (Argyle and Dean, 1965) proposes that interagent spacing is governed by an equilibrium of intimacy between two people, consisting of the amount of mutual eye gaze and smiling, the intimacy of topic or thought, and posture and arm configurations. Interpersonal Adaptation Theory (Burgoon et al., 1995) extends this by considering more dynamic behaviors—such as the amount of convergence, divergence, mirroring, compensation, and reciprocity—between social partners. Computational models based on psychological parameters have been investigated in immersive virtual environments (Bailenson et al., 2001) and HRI (Mumm and Mutlu, 2011, Takayama and Pantofaru, 2009).

There is little work bridging the gap between physical and psychological representations of proxemics. The psychological relationship between two social agents dictates the desired sensory experience between them, manifested physically through change in position and orientation (Figure 2.1). This dynamic can be expressed using a psychophysical representation, situating the agents in the interaction and the environment (Mead and Matarić, 2012, Mead et al., 2012).

2.3.3.3 Psychophysical Representation of Proxemic Behavior

The *psychophysical* representation is concerned with the perception and production of social stimuli by two or more interacting agents. The four proxemic zones proposed by Hall (1966)—public, social, personal, and intimate—describe the interpersonal relationship between two people and, thus, utilize a *psychological* representation of proxemics. However, Hall (1963, 1966) posited that these zones are characterized by *psychophysical* factors, such as the visual, auditory (voice loudness), olfactory, thermal, and somatosensory (touch and kinesthetic) experiences of each interacting participant.

For example, upon first meeting, two Western American strangers often shake hands, and, in doing so, subconsciously estimate each other’s arm length; these strangers will then stand just outside of the extended arm’s reach of the other, so as to maintain a safe distance from a potential fist strike (Hediger, 1955). This sensory experience characterizes “social distance” between strangers or acquaintances. As their relationship develops into a friendship, the risk of a fist strike is reduced, and they are willing to stand within an arm’s reach of one another at a “personal distance”; this is highlighted by the fact that brief physical embrace (e.g., hugging) is common at this range (Hall, 1963). However, olfactory and thermal sensations of each other are often not as desirable in a friendship, so some distance is still maintained to reduce the potential of these sensory experiences. For these sensory stimuli to become more desirable and

social, the relationship would have to become more intimate; olfactory, thermal, and prolonged tactile interactions are characteristic of intimate interactions, and can only be experienced at close range, or “intimate distance” (Hall, 1963).

The psychophysical representation relates sensory experience of multimodal social stimuli (e.g., speech and gesture) to physical parameters (e.g., distance and orientation), and also serves as a bridge between the physical and psychological proxemic representations, situating sociable agents (both human and robot) in the interaction and the environment (Figure 2.1). However, the psychophysical representation has not yet been thoroughly examined, formalized, and adopted in HRI. This dissertation investigates the use of the psychophysical representation for unifying proxemic behavior and multimodal communication in autonomous face-to-face human-robot interactions.

2.4 Summary

This chapter surveyed the field of human-robot interaction with a focus on autonomous sociable and socially assistive robotics (Sections 2.2.1 and 2.2.2, respectively). A review of proxemic behavior in both human-human and human-robot interactions was provided (Sections 2.3.1 and 2.3.2, respectively), and features commonly used in proxemic behavior analysis were consolidated into three representations of proxemic behavior: physical, psychological, and psychophysical (Section 2.3.3). The psychophysical representation encodes the sensory experience of social stimuli (e.g., speech and gestures) by each agent engaged in a face-to-face social interaction. This representation motivated the formalization of a computational framework unifying proxemic behavior and multimodal communication in human-robot interaction, described in the next chapter.

Chapter 3

Framework for Proxemics and Multimodal Communication

This chapter describes the approach for a novel representation of proxemics and multimodal communication, and its formalization in a computational framework for human-human and human-robot interactions.

3.1 Overview

The main contribution of this dissertation is in the formalization of a novel representation that unifies proxemics and multimodal communication, and its integration into an extensible computational framework for human-human and human-robot interactions.

To achieve these goals, two approaches are formalized: a *heuristic approach* and a *data-driven approach*. The heuristic approach establishes computational representations of features that come from related work in human-human proxemic behavior analysis, including psychophysical considerations that estimate how humans perceive multimodal communication in different proxemic configurations. Insights gained from the heuristic approach motivated the data-driven approach to human-robot proxemics.

Heuristic Approach: A heuristic approach based on proxemic representations derived from social science literature is presented in Section 3.2. This framework first considers the body pose (i.e., the position and orientation of the head, shoulders, torso, hip, and stance) of each individual sociable agent (human and robot). Individual poses of interacting agents are then used to extract physical distance and orientation between each social dyad (pair of agents). These physical distance and orientation parameters are used to estimate the psychophysical (visual, voice loudness, kinesthetic, olfactory, thermal, and touch) experiences of each agent based on a coding schema from the social sciences (Hall, 1963). Finally, the psychophysical features are used to determine a categorical representation of the interpersonal relationship between two agents (Hall, 1963, 1966). The heuristic framework is representative of human-human proxemic preferences (Section 2.3.1), and extends approaches taken by related work in HRI (Section 2.3.2).

Data-Driven Approach: A data-driven approach is presented in Section 3.3. This framework establishes proxemic preferences for the robot (as opposed to the human) based on how it would perceive human social signals (speech and gestures) in HRI. The framework is motivated by insights gained from the heuristic approach, is inspired by the psychophysical processes that govern human-human proxemic behavior (Hall, 1963, 1966), and is designed to improve robot performance (measured by automated social signal recognition rates) in HRI. The framework consists of three primary components:

1. the proxemic (distance and orientation) preferences of all agents (human and robot) with respect to other agents,
2. predictive models of how social signals will be *produced* by each transmitting agent in different proxemic configurations, and
3. predictive models of how social signals will subsequently be *perceived* by each receiving agent.

3.2 Heuristic Approach

In the heuristic approach, proxemic behavior is represented by a set of “features” based on the most commonly used metrics in the social sciences literature. Features in Schegloff’s *individual* representation (Schegloff, 1998) are extracted first, which represent the body configuration of each social agent (human or robot). These individual features are used to calculate the features for each social dyad (pair of social agents). This dissertation focuses on three popular and validated annotation schemas for proxemic behavior analysis: 1) Mehrabian’s *physical* feature representation (Mehrabian, 1972), 2) Hall’s *psychophysical* feature representation (Hall, 1963), and 3) Hall’s *psychological* feature representation (Hall, 1966).

3.2.1 Features of Individual Representation

The *individual* representation extracts the body poses of each interacting agent. While in itself not a representation of proxemics, Schegloff (1998) emphasized the importance of distinguishing between relative poses of the lower and upper parts of the body of an individual during face-to-face social interactions (Figure 3.1), suggesting that changes in the lower parts (from the waist down) signal *dominant involvement*, while changes in the upper parts (from the waist up) signal *subordinate involvement*. When a pose deviates from its *home position* (i.e., 0°) with respect to an “adjacent” pose, the deviation does not last long and a compensatory orientation behavior occurs, either from the subordinate or the dominant body part. More often, the subordinate body part (e.g., head) is responsible for the deviation and, thus, provides the compensatory behavior; however, if the dominant body part (e.g., shoulder) is responsible for the deviation or provides the compensatory behavior, a shift in attention (or *involvement*) is likely to have occurred. Schegloff (1998) referred to this phenomenon as *body torque*, which has been investigated in HRI (Kuzuoka et al., 2010).

The following individual features were represented and used in this work:

- **Stance Pose:** most dominant involvement cue; position midway between the left and right ankle positions and orientation orthogonal to the line segment connecting the left and right ankle positions
- **Hip Pose:** subordinate to stance pose; position midway between the left and right hip positions, and orientation orthogonal to the line segment connecting the left and right hip positions
- **Torso Pose:** subordinate to hip pose; position of torso and average of hip pose orientation and shoulder pose orientation (weighted based on relative torso position between hip pose and shoulder pose)
- **Shoulder Pose:** subordinate to torso pose; position midway between the left and right shoulder positions and orientation orthogonal to the line segment connecting the left and right shoulder positions
- **Head Pose:** subordinate to shoulder pose; can be extracted and tracked using head pose estimation techniques, such as that described in (Morency et al., 2008)

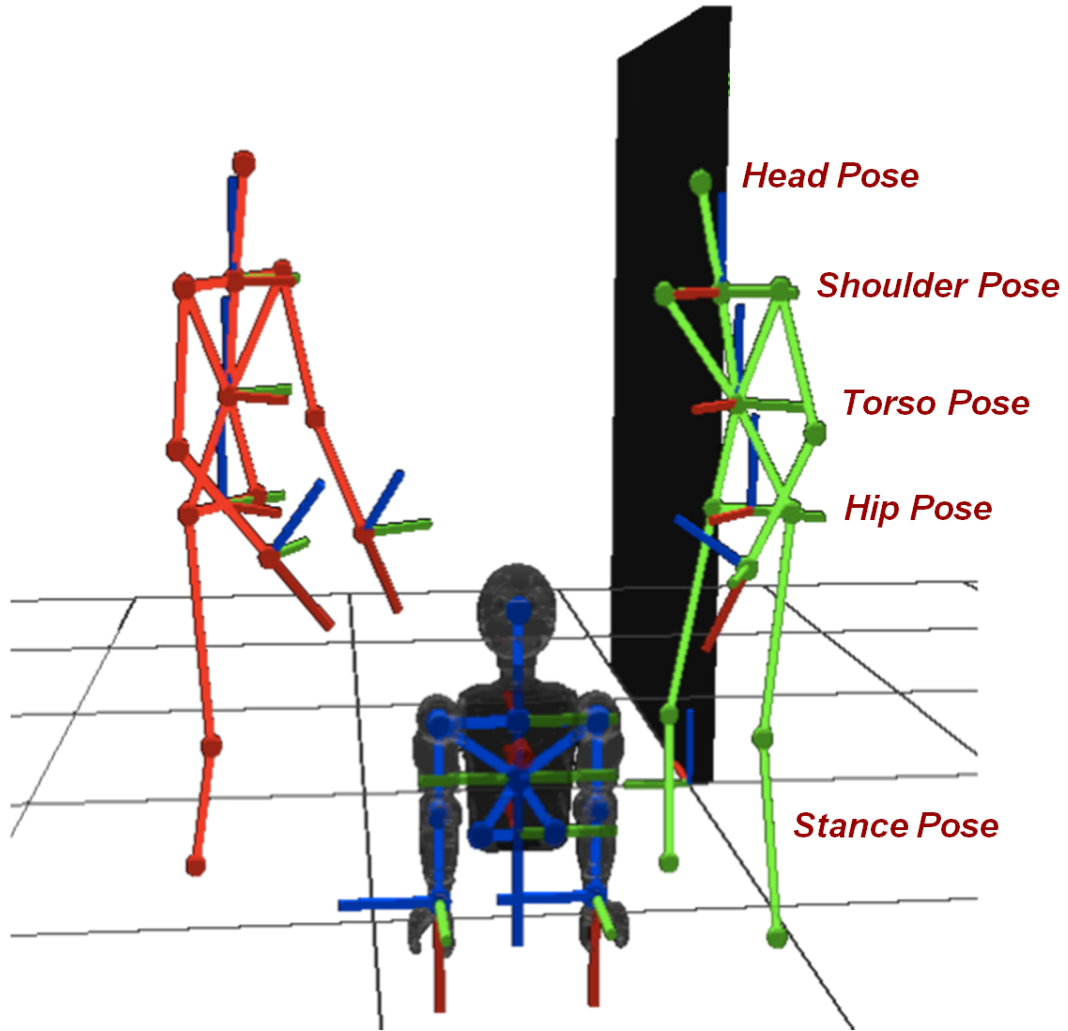


Figure 3.1: Individual pose features for two human users and an upper-body humanoid robot within a single framework; the absence of some features—such as the head, arms, or legs—signifies a pose estimate with low confidence.

3.2.2 Features of Physical Representation

Mehrabian (1972) provides distance- and orientation-based metrics between a social dyad (two individuals) for proxemic behavior analysis (Figure 3.2). These *physical* features are the most commonly used in the study of both human-human and human-robot proxemics.

The following annotations are made for each individual in a dyad between social agents A and B (human or robot):

- **Total Distance:** magnitude of a *Euclidean distance vector* from the pelvis of agent A to the pelvis of agent B
- **Straight-Ahead Distance:** magnitude of the x -component of the total distance vector
- **Lateral Distance:** magnitude of the y -component of the total distance vector
- **Relative Body Orientation:** magnitude of the angle of the pelvis of agent B with respect to the pelvis of agent A

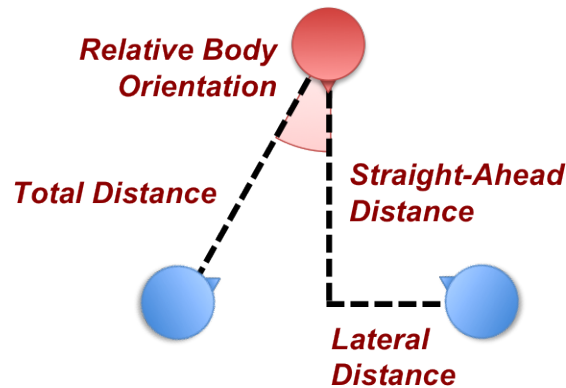


Figure 3.2: In this interaction scenario, proxemic behavior is analyzed using simple physical features between each social dyad (pair of individuals).

3.2.3 Features of Psychological Representation

Hall (1966) provides a *psychological* proxemic representation that uses physical distance and orientation features to categorize the interpersonal relationship between people in a social dyad into two proxemic codes: *distance* and *sociofugal-sociopetal (SFP) axis* (Figure 3.3).

The psychological “feature *codes*” and “feature *intervals*” are annotated as follows:

- **Distance Code¹:** based on total distance; intimate (0”–18”), personal (18”–48”), social (48”–144”), public (more than 144”)
- **Sociofugal-Sociopetal (SFP) Axis Code:** based on relative body orientation (in 20° intervals), with face-to-face (axis-0) representing maximum sociopetality and back-to-face (axis-8) representing maximum sociofugality (Lawson, 2001, Low and Lawrence-Ziga, 2003, Sommer, 1967); axis-0 (0°–20°) axis-1 (20°–40°), axis-2 (40°–60°), axis-3 (60°–80°), axis-4 (80°–100°), axis-5 (100°–120°), axis-6 (120°–140°), axis-7 (140°–160°), or axis-8 (160°–180°)

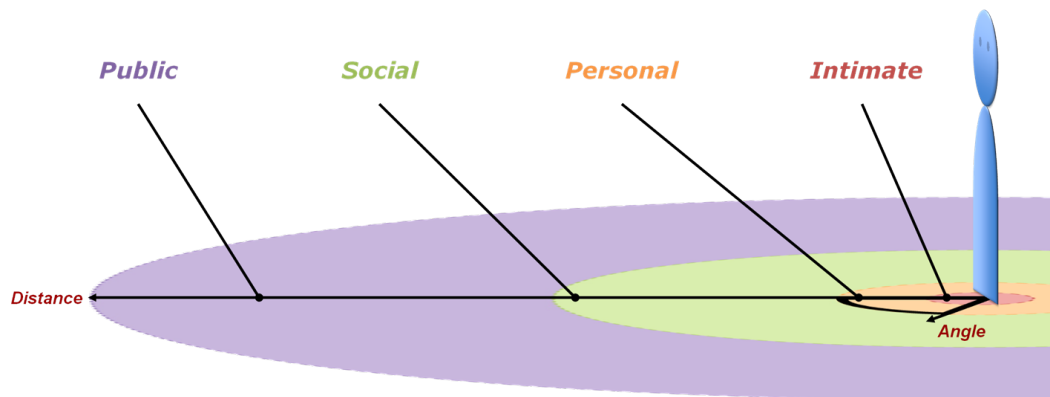


Figure 3.3: Public, social, personal, and intimate distance codes, and SFP axis codes.

¹These proxemic distances pertain to Western culture—they are not cross-cultural (Hall, 1963, 1966).

3.2.4 Features of Psychophysical Representation

Hall (1963) proposes a *psychophysical* proxemic representation as an alternative to strict physical analysis and is more descriptive than the psychological representation, providing a sort of functional sensory explanation to the human use of space in multimodal social interaction (Figure 3.4). Hall (1963) seeks not only to answer questions of *where* a person will be, but, also, the question of *why* they are there, addressing the underlying processes and systems that govern proxemic behavior. Examples of the utility of this representation are provided in Section 2.3.3.3. This dissertation utilizes and extends this approach for robots socially interacting with human users.

The coding schema proposed by Hall (1963) is typically annotated by social scientists based purely on distance and orientation data observed from video (Hall, 1966). The automation of this tedious process is a major contribution of the dissertation.

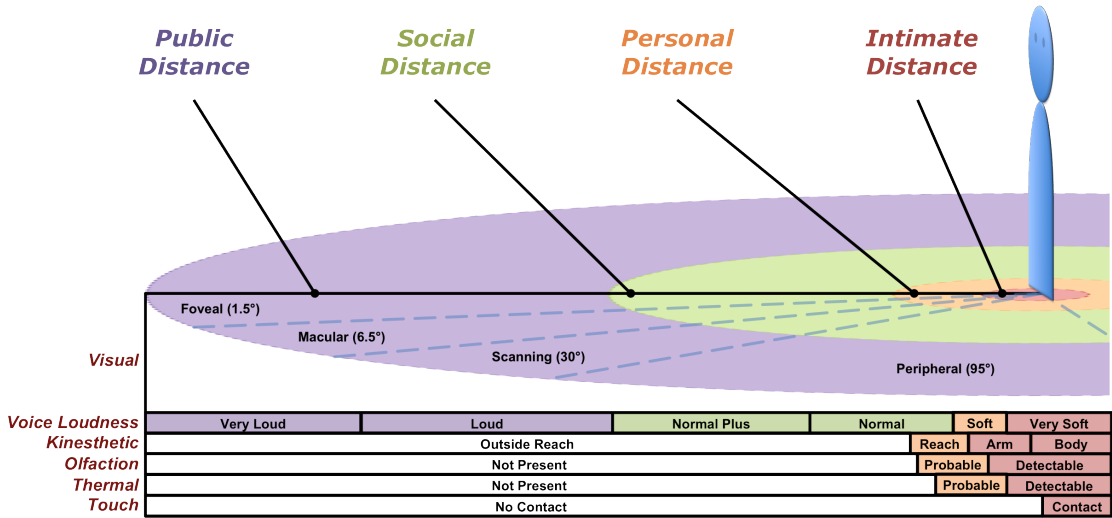


Figure 3.4: The anticipated sensory sensations that an individual would likely experience in different physical proxemic configurations and within the psychological distance zones.

The psychophysical “feature *codes*” and the corresponding “feature *intervals*” for each individual in a dyad between social agents *A* and *B* (human or robot) are as follows:

- **Visual Code:** based on head pose; foveal (sharp; 1.5° off-center), macular (clear; 6.5° off-center), scanning (30° off-center), peripheral (95° off-center), or no visual contact
- **Voice Loudness Code:** based on total distance; silent (0”–6”), very soft (6”–12”), soft (12”–30”), normal (30”–78”), normal plus (78”–144”), loud (144”–228”), or very loud (more than 228”)
- **Kinesthetic Code:** based on the distances between the hip, torso, shoulder, and head poses, and measured arm length; within body contact distance, within easy touching distance with only forearm extended, within reaching distance, or outside reaching distance
- **Olfaction Code:** based on total distance; body odor or breath detectable (0”–18”), olfaction probably present (18”–36”), or olfaction not present
- **Thermal Code:** based on total distance; conducted or radiant heat detected (0”–12”), heat probably detected (12”–21”), or heat not detected
- **Touch Code:** based on total distance; contact or no contact

The psychophysical representation is the crucial component in the unification of proxemics and multimodal communication. It offers insights into how social signals produced—either intentionally (e.g., speech or gesture) or unintentionally (e.g., smell or body heat)—by a sociable agent will likely be perceived (with respect to *intensity*, not *interpretation*) by other interacting sociable agents in any proxemic configuration. This is the fundamental philosophy adopted by this dissertation.

3.2.5 Discussion

The heuristic approach establishes a consolidated computational framework for proxemic behavior analysis based on individual (Schegloff, 1998), physical (Mehrabian, 1972), psychological (Hall, 1966), and psychophysical (Hall, 1963) metrics from the social sciences. Most notably, the formalization of the psychophysical representation unifies proxemic behavior and social stimuli present in face-to-face multimodal social interactions. However, the values presented in the psychophysical representation only represent human perceptions of social stimuli; while this approach might suffice for explaining human-human interactions and could be used to approximate human-robot interactions, robots will not necessarily perceive those stimuli in the same way—an alternative approach is needed to determine how robots will perceive social signals, and subsequently how it will impact robot proxemic behavior in HRI. A data-driven approach to resolve this issue is presented in the next section.

3.3 Data-Driven Approach

In this dissertation, the data-driven approach contributes a novel probabilistic proxemic framework, motivated by the psychophysical representation described in Section 3.2.4, that considers how all represented sociable agents—both humans and robots—experience a co-present interaction (Mead and Matarić, 2012). Models of the production (output) and perception (input) of speech and gesture are conditioned on interagent pose (position and orientation). A definition of the framework parameters is provided, followed by a model representation relating these parameters to one another.

3.3.1 Definition of Framework Parameters

Consider two sociable agents, A and B , that are co-located and intend to interact. At any point in time and from any location in the environment, an agent must be capable of estimating:

1. An **interagent pose**, POS —Where will B stand relative to A ?
2. A **speech output level**, SOL_{BA} —How loudly will B speak to A ?
3. A **gesture output level**, GOL_{BA} —In what space will B gesture to A ?
4. A **speech input level**, SIL_{AB} —How well will A perceive B 's speech?
5. A **gesture input level**, GIL_{AB} —How well will A perceive B '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—they are not about *what* is said, but, rather, *how* it is said.

Interagent pose (POS) is expressed as a toe-to-toe distance (d) and two orientations—one from A to B (α), and one from B to A (β). For navigation and trajectory planning, one typically works in world coordinates (Marder-Eppstein et al., 2010); however, for proxemics, it is beneficial to consider agent poses relative to each another (the coordinate transforms are computationally inexpensive).

Speech output and **input levels** (SOL_{BA} and SIL_{AB}) are each represented as a sound pressure level, a logarithmic measure of sound pressure relative to a reference value (ambient noise in the environment), 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 SOL_{BA} and SIL_{AB} .

Gesture output and **input levels** (GOL_{BA} and GIL_{AB}) are each represented as a 3D region of space called a gesture locus (McNeill, 1992, Rossini, 2004). For gesture production (gesture output levels), the locus defines the regions of the body, along Cartesian axes, involved by the gesture; Rossini (Rossini, 2004) discretized these values along the coronal (x), sagittal (y), and transverse (z) body axes. In this dissertation, the GOL_{BA} locus is modeled as a continuous distribution with respect to an agent frame (e.g., a distribution of body part locations with respect to the agent’s base pose). Related work in HRI suggests that robot nonverbal behaviors (including gestures) should be parameterized based on proxemics (Brooks and Arkin, 2007). For gesture production, the GOL_{BA} locus is modeled as the locations of B ’s body parts in physical space. The GIL_{AB} is then be modeled as the region of A ’s visual field occupied by the body parts associated with the gesture output of B (i.e., GOL_{BA}) (Mead and Matarić, 2012).

3.3.2 Modeling Framework Parameters

In this work, distributions of these pose, speech, and gesture parameters and their relationships are modeled as a Bayesian network to represent (Figure 3.5):

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.

Formally, each component of the model can be written respectively as:

$$p(POS) \tag{3.1}$$

$$p(SOL_{BA}, GOL_{BA} | POS) \tag{3.2}$$

$$p(SIL_{AB}, GIL_{AB} | SOL_{BA}, GOL_{BA}, POS) \tag{3.3}$$

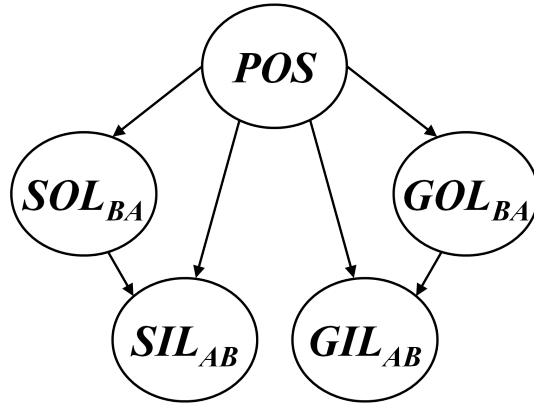


Figure 3.5: Bayesian network modeling relationships between pose, speech, and gestures.

3.3.3 Discussion

The data-driven approach is inspired by the relationship between the physical representation (describing a proxemic configuration; Section 3.2.2) and the psychophysical representation (how multimodal communication is produced and perceived in different proxemic configurations; Section 3.2.4) first formalized in the heuristic approach (Section 3.2). The heuristic approach was representative of the sensory sensations of individuals during human-human interactions, and could be used as an approximation for human-robot interactions. The data-driven approach provides a computational framework in which models of robot social signal perception can be estimated based on predictions of how a human user will produce social signals in different proxemic configurations, thus serving as a more situated alternative to the heuristic approach.

3.4 Summary

This chapter discussed a unifying framework for situated proxemics and multimodal communication in face-to-face human-robot interactions. Two approaches—heuristic and data-driven—were presented, each formalizing traditional representations of proxemics and communication into a single unifying computational framework. The heuristic approach is derived from social science literature on human-human interactions; while this approach could be used to approximate robot behavior in human-robot interactions, it serves primarily to provide inspiration and guidelines for the derivation of computational models for robot social signal perception and proxemics in the data-driven approach. The next chapter discusses the procedure for extracting features and modeling parameters in the heuristic and data-driven approaches, respectively.

Chapter 4

Modeling Proxemics and Multimodal Communication

This chapter discusses the process of modeling parameters of the proxemic and multimodal communication framework described in the previous chapter. The heuristic approach is modeled based on values from literature and the sensor selected for feature detection; the approach enables the automated extraction of commonly used discrete proxemic features. The data-driven approach is modeled based on data collected in both human-human and human-robot interactions, and serves as a continuous alternative to the heuristic approach.

4.1 Overview

In the previous chapter, a computational framework for unifying proxemic behavior and multimodal communication was discussed. The framework highlights a psychophysical representation, estimating how social stimuli (present in multimodal communication) are perceived by sociable agents (human and robot) in different proxemic configurations (distances and orientations). Two approaches—heuristic and data-driven—were presented to integrate psychophysical considerations into models for HRI.

The next sections discuss the processes by which both the heuristic and data-driven approaches are modeled in the framework. For the heuristic approach, error models of extracted features are discussed. For the data-driven approach, a data collection and analysis is provided to inform probabilistic graphical models of proxemic behavior and multimodal communication. The resulting models of these approaches have implications for autonomous human behavior recognition (implemented and evaluated in Section 5.2) and autonomous robot behavior control (implemented and evaluated in Section 5.3).

4.2 Models of the Heuristic Approach

The heuristic approach of proxemic feature extraction can be implemented using any human motion capture technique; this work utilized the PrimeSensor structured light hardware range sensor and the OpenNI¹ person tracking software for markerless motion capture. This setup was chosen because 1) it is non-invasive to participants, 2) it is readily deployable in a variety of environments (ranging from an instrumented workspace to a mobile robot), and 3) it does not interfere with the interaction itself. Joint pose estimates provided by this setup were used to extract individual features, which were then used to extract the physical and psychophysical features of each interaction dyad (two individuals). Error models were developed based on data collections (the data for which are publicly available upon request) that evaluate the automated extraction of the individual (Schegloff, 1998), physical (Mehrabian, 1972), psychological (Hall, 1966), and psychophysical (Hall, 1963) feature representations, discussed below.

¹<http://wiki.ros.org/openni>

4.2.1 Modeling Features of Individual Representation

The precision of the PrimeSensor distance estimates was modeled and analyzed. The PrimeSensor was mounted atop a tripod at a height of 1.5 meters and pointed straight at a wall. The sensor rig was placed at 0.2-meter intervals between 0.5 meters and 2.5 meters and distance readings (a collection of 3-dimensional points, referred to as a “point cloud”) were taken at each location. A planar model segmentation technique² was used to eliminate points in the point cloud that did not fit onto the wall plane. The average depth reading of each point in the segmented plane was calculated, and the sensor error E as a function of distance d (in meters) was modeled as $E(d) = k \times d^2$, with $k = 0.0055$. The procedure and results are consistent with those reported in accuracy analyses of other structured light sensors³, which often differ only in the value of k ; thus, if a similar range sensor were to be used, system performance would scale with k .

Shotton et al. (2011) provides a comprehensive evaluation of individual joint pose (e.g., head, shoulder, torso, etc.) estimates produced by the Microsoft Kinect, a structured light range sensor very similar in hardware to the PrimeSensor. The study reports an accuracy of 0.1 meters for each detected human joint, and a mean average precision of 0.984 and 0.914 for tracked (observed) and inferred (obstructed or unobserved) joint estimates, respectively. While the underlying algorithms differ, the performance was comparable in this dissertation work.

²http://www.pointclouds.org/documentation/tutorials/planar_segmentation.php

³http://www.ros.org/wiki/openni_kinect/kinect_accuracy

4.2.2 Modeling Features of Physical Representation

In the course of estimating subsequent dyadic proxemic features, it is important to note that any range sensor detects the *surface* of the individual and, thus, joint pose estimates are projected into the body by some offset. In Shotton et al. (2011), this value is determined from data, with an average offset of 0.039 meters. To extract accurate physical proxemic features of the social dyad, twice this value (once for each individual) was subtracted from the measured ranges to determine the surface-to-surface distance between two bodies. A comprehensive data collection and analysis of the joint pose offset used by the system is reported in Shotton et al. (2011).

4.2.3 Modeling Features of Psychological Representation

Intervals for each psychological feature code were evaluated at 1, 2, 3, and 4 meters from the sensor; the sensor was orthogonal to a line passing through the hip poses of two standing individuals. The results of these estimates are illustrated in Figures 4.1 and 4.8. At some ranges from the sensor, a feature interval places the individuals out of the sensor field-of-view; these values are omitted⁴.

The evaluations of feature intervals for the *distance code* (Figure 4.1) demonstrate the effect of interval size on feature classification—the larger “social” interval was extracted with higher precision (i.e., there is more room for error in the sensor reading).

The *SFP axis code* contains feature intervals of uniform size (40°). Rather than evaluate each interval independently, the average precision of the intervals was evaluated at different distances. The error in estimated orientation of one individual with respect to another individual was considered at 1, 2, 3, and 4 meters from each other. At shorter ranges, the error in estimated position of one individual increases the uncertainty of the orientation estimate (Figure 4.2).

⁴This occurs for the SFP axis estimates, as well as at the public distance interval.

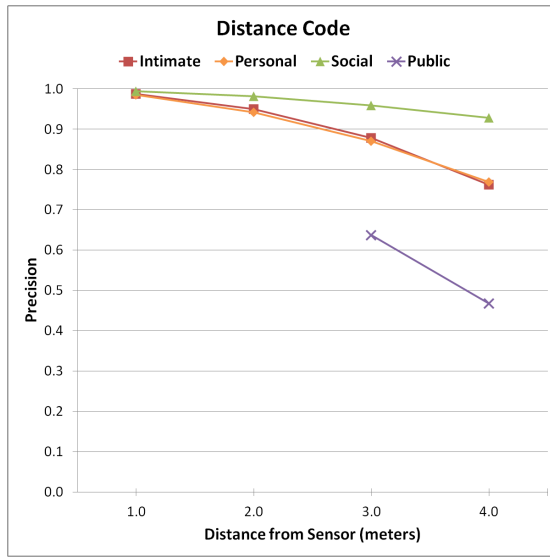


Figure 4.1: Error model of psychological distance code.

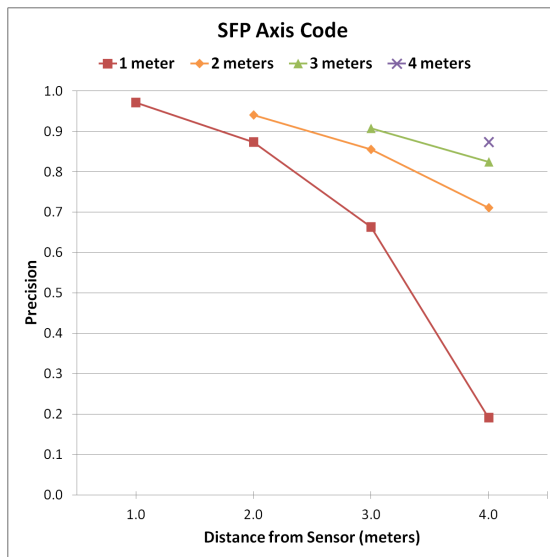


Figure 4.2: Error model of psychological sociofugal-sociopetal (SFP) axis code.

4.2.4 Modeling Features of Psychophysical Representation

Each feature annotation in Hall’s psychophysical representation (Hall, 1963) was developed based on values from literature on the human sensory system (Hall, 1966). It is beyond the scope of this work to evaluate whether or not a participant actually experiences the stimulus in the way specified by a particular feature interval⁵—such evaluations come from literature cited by Hall (Hall, 1963, 1966) when the representation was initially proposed. Rather, this dissertation provides a *theoretical* error model of the psychophysical feature annotations as a function of their respective distance and orientation intervals based on the sensor error models provided above.

As with the psychological representation, intervals for each feature code were evaluated at 1, 2, 3, and 4 meters from the sensor; the sensor was orthogonal to a line passing through the hip poses of two standing individuals. The results of these estimates are illustrated in Figures 4.3–4.8. At some ranges from the sensor, a feature interval places the individuals out of the sensor field-of-view; these values are omitted⁶.

The evaluations of feature intervals for the *voice loudness* (Figure 4.4), *olfaction* (Figure 4.6), *thermal* (Figure 4.7), and *touch* (Figure 4.8) codes demonstrate the effect of interval size on feature classification. Specifically, the larger and/or more common feature intervals are extracted with higher precision, suggesting their strong utility for annotating proxemics in typical social encounters.

⁵For example, the radiant heat or odor transmitted by one individual and the intensity at the corresponding sensory organ of the receiving individual is not measured.

⁶This occurs for the loud and very loud voice loudness intervals, as well as for the outside reach kinesthetic interval.

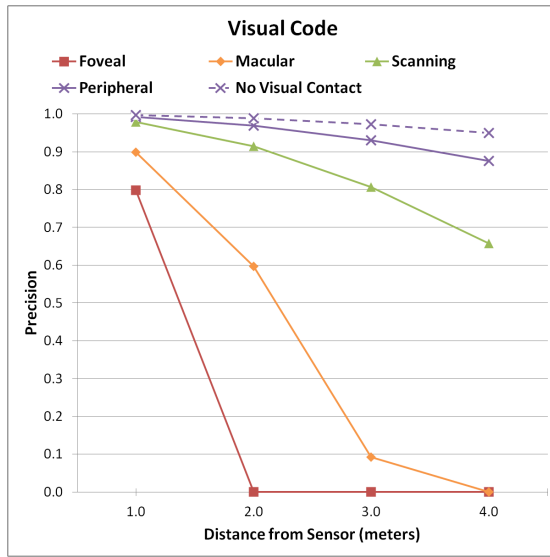


Figure 4.3: Error model of psychophysical visual code.

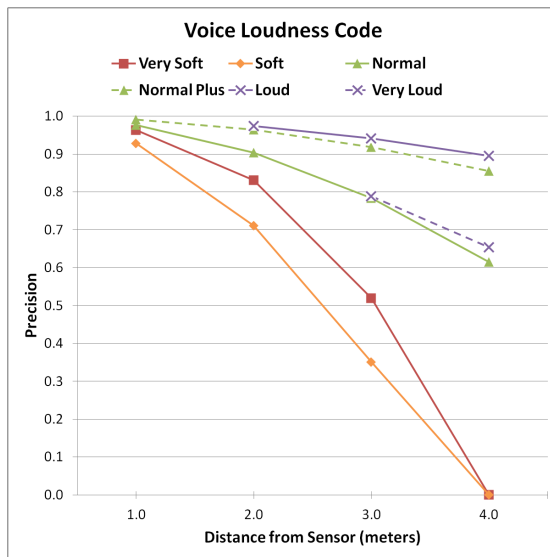


Figure 4.4: Error model of psychophysical voice loudness code.

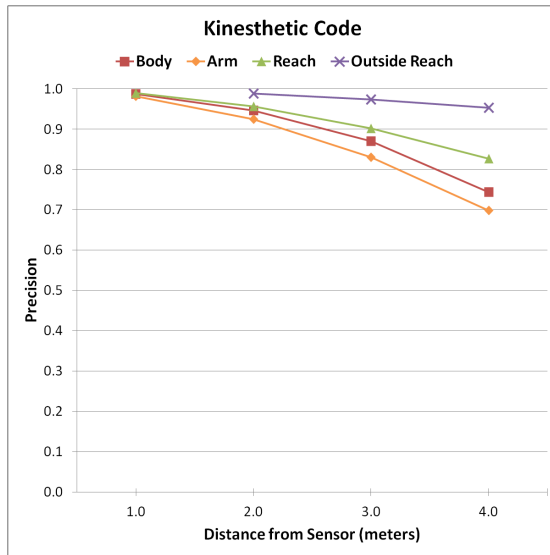


Figure 4.5: Error model of psychophysical kinesthetic code.

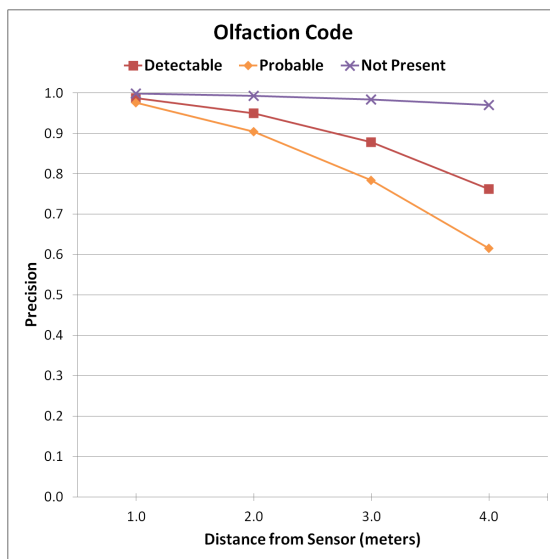


Figure 4.6: Error model of psychophysical olfaction code.

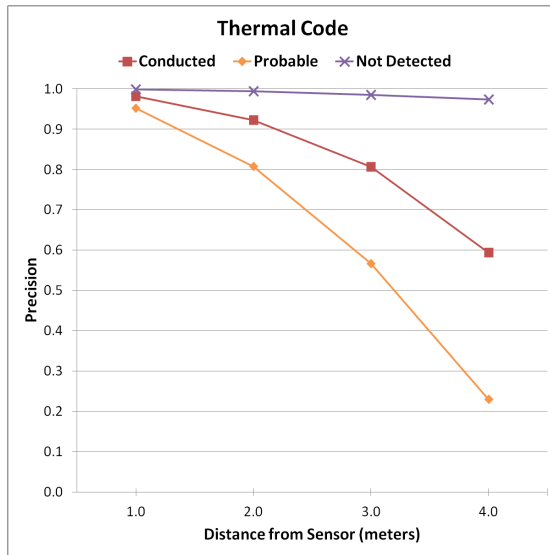


Figure 4.7: Error model of psychophysical thermal code.

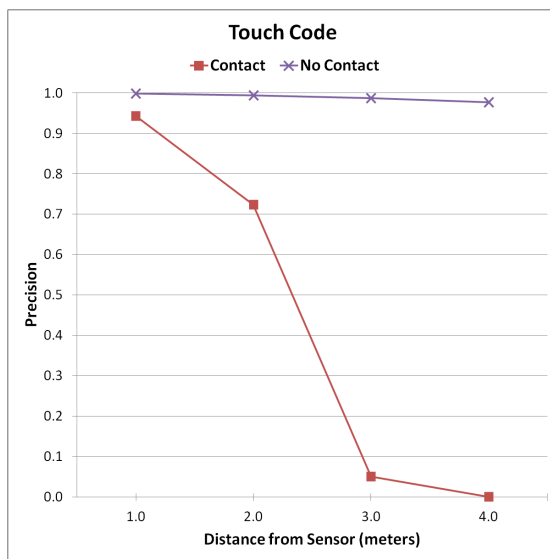


Figure 4.8: Error model of psychophysical touch code.

Eye gaze was unavailable for true *visual code* extraction. Instead, the *visual code* was estimated based on the head pose estimation (Morency et al., 2008); when this approach failed, shoulder pose was used to estimate eye gaze. Both of these estimators resulted in coarse estimates of narrow feature intervals (foveal and macular) (Figure 4.3).

For evaluation of the *kinesthetic code*, interval distance estimates based on average human limb lengths (calculated based on readings from the PrimeSensor) were used (Figure 4.5) (Hall, 1966). In practice, performance is expected to be lower, as the feature intervals are variable—they are calculated dynamically based on joint pose estimates of the individual (e.g., the hips, torso, neck, shoulders, elbows, and hands). The uncertainty in joint pose estimates accumulates in the calculation of the feature interval range, so there are expected to be more misclassifications at the feature interval boundaries.

4.2.5 Discussion

The models used in this heuristic approach are based on empirical measures provided by the social science literature (Hall, 1963, 1966, Mehrabian, 1972, Schegloff, 1998), resulting in a coarse discretization of the parameter space in the psychological (Hall, 1966) and psychophysical (Hall, 1963) representations. While this might be adequate for human-human proxemic behavior analysis (and human proxemic behavior recognition, discussed in Section 5.2), the approach does not address the problem at the fine-grained resolution necessary for control of robot proxemic behavior and multimodal communication in human-robot interactions (discussed in Section 5.3). The next section presents a more continuous approach for HRI inspired by the psychophysical representation (specifically, speech and gesture production and perception) and derived from data.

4.3 Models of the Data-Driven Approach

A data collection was conducted and data were analyzed to inform the parameters of the models in Equations 3.1 – 3.3. The procedure was performed in the context of face-to-face human-human and human-robot interactions at controlled distances. In the data collection and analysis, a human and robot agents are denoted H and R , respectively.

4.3.1 Data Collection

4.3.1.1 Procedure

Each participant watched a short (1–2 minute) cartoon at separate stations (denoted $C1$ and $C2$ in Figure 4.9), 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 gesture frequently when describing them (McNeill, 1992) (Figure 4.10). Participants interacted for 4–6 minutes, 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 $p(POS)$ (Equation 3.1), and also informing speech and gesture output/input parameters conditioned upon those natural poses (Equations 3.2 and 3.3). 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 Figure 4.9), 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 models represented in Equations 3.2 and 3.3.

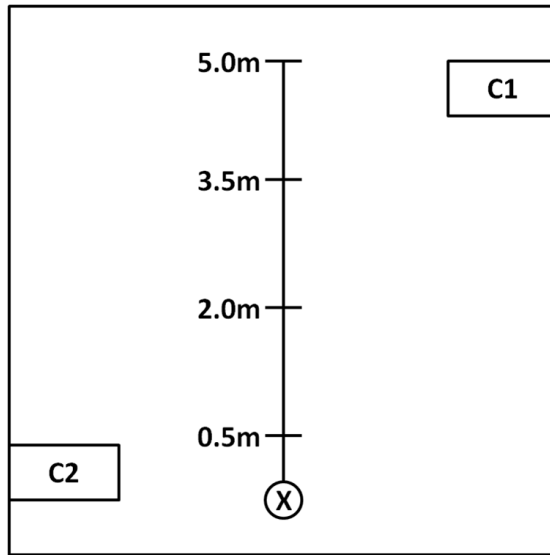


Figure 4.9: The experimental setup for modeling proxemic behavior and multimodal communication. Participants watched cartoons at locations $C1$ and $C2$. After watching each cartoon, one social partner (either human or robot; **human vs. robot condition**) relocated to the floor mark X . A participant approached the social partner along the line, and either a) stopped in at any interagent distance (**natural distance conditions**), or b) at one of four specified distances ($d = \{0.5, 2.0, 3.5, 5.0\}$ meters; **controlled distance conditions**).

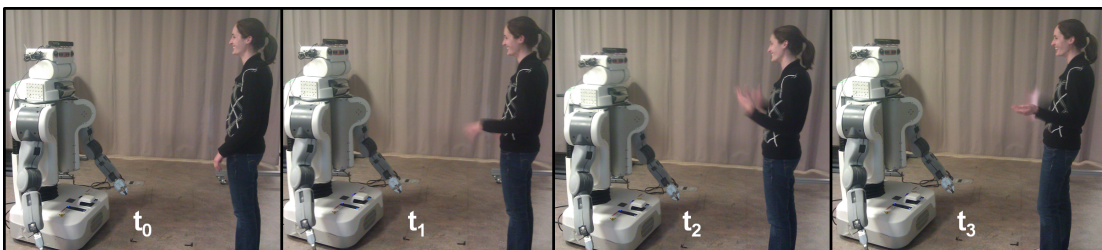


Figure 4.10: A participant uses gestures while describing a cartoon to the PR2 robot.

For each participant, the “social partner” was either another human participant or a PR2 robot (“**human vs. robot**” **between-participants condition**) (robot shown in Figure 4.10). For each interaction, participants watched and discussed three different cartoons and three of the same cartoons to explore the space of multimodal 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.

4.3.1.2 Materials and Measures

Five Microsoft Kinects strategically positioned along the line of approach (Figure 4.9) were used to monitor interagent pose (POS), as well as to extract the positions of human body parts (e.g., head, torso, and hands), which provided the representation of human gesture output levels (GOL_{HR}) (Mead et al., 2013). 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.) (Mead and Matarić, 2012) (Figure 4.11). 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.

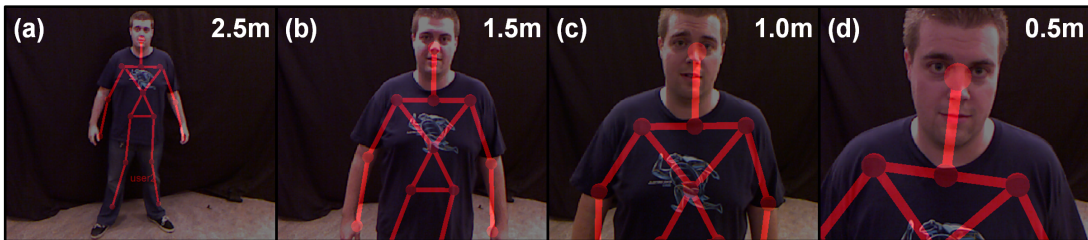


Figure 4.11: Body features that fall into the Kinect field-of-view, depicted at four distances: (a) 2.5 meters, (b) 1.5m, (c) 1.0m, and (d) 0.5m (Mead and Matarić, 2012).

4.3.1.3 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 in a dyad 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 the goal was to model framework parameters over a general population.

4.3.1.4 Dataset

Pose, speech, and gesture data were collected for 10 human-human interactions (20 participants) and 20 human-robot interactions (20 participants); these data are publicly available upon request. 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 sections.

4.3.2 Data Modeling and Analysis

The resulting dataset was used to model the relationships between interagent pose (Equation 3.1) and speech/gesture output and input levels (Equations 3.2 and 3.3).

4.3.2.1 Proxemic Preferences: Human Interagent Pose

To inform Equation 3.1, 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** (Figure 4.12).

The interagent poses in the human-human interactions were consistent with human-human proxemics literature—participants positioned themselves at a “social distance” (1.2–3.7 meters), as predicted by their interpersonal relationship (strangers or, at most, acquaintances) (Hall, 1963).

However, the interagent poses in the human-robot interactions were inconsistent with human-robot proxemics literature—participants in the data collection positioned themselves much farther away than has been observed in related work. Walters et al. (2009) 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 and Pantofaru (2009) investigated spacing between humans and the PR2 (the same robot that was used in this dissertation study), reporting mean distances of 0.25–0.52 meters. Differences between these 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 the dissertation study, participants were more focused on the interaction itself, so the positioning might have been more natural. Second, in many of these studies, the robot is not producing gestures; however, in this dissertation study, the robot was gesturing and had a long reach (0.92 meters), so participants might have positioned themselves farther away from the robot to avoid physical contact.

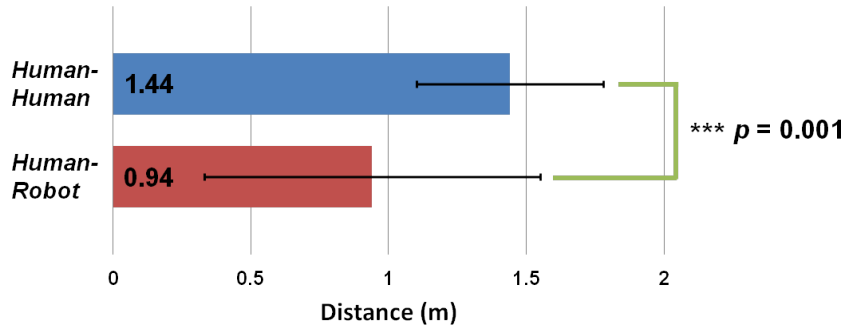


Figure 4.12: Comparisons of human-human vs. human-robot proxemics ($p = 0.001$).

4.3.2.2 Social Signal Production: Human Speech/Gesture Output Levels

To inform Equation 3.2, human speech and gesture output levels (SOL_{HR} and GOL_{HR} , respectively) were considered. 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, a trend towards significance ($p = 0.073$) was found for the statistical interaction between distance and **human vs. robot conditions** (Figure 4.13) 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 Kriz et al. (2010), 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; the dissertation work expands upon this result by illustrating the effect at multiple distances.

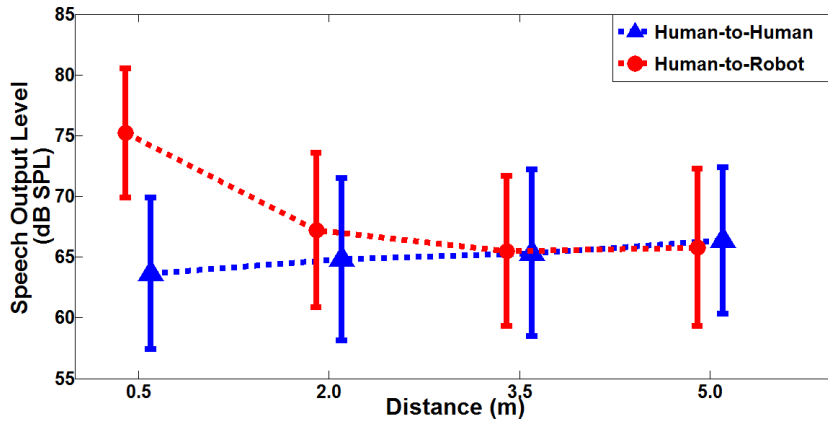


Figure 4.13: Human speech output levels vary with distance in human vs. robot conditions.

No significant differences in gesture output levels were detected 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 this 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 (Kendon, 1990)) influences the way in which people produce gestures (Ozyurek, 2002). For example, during a face-to-face interaction, one might expect to see gestures made directly in front of a person (Ozyurek, 2002), likely in the center of the field-of-view of the robot; however, during a non-frontal interaction, one might expect to see gestures located more laterally (Ozyurek, 2002), potentially falling out of the field-of-view of the robot. This consideration informs the selection of hardware sensors for vision-based human gesture recognition in HRI—cameras with a wide horizontal field-of-view are desirable.

4.3.2.3 Social Signal Perception: Robot Speech/Gesture Input Levels

To inform Equation 3.3, data recorded by the microphones and camera of the Kinect on-board the robot were used to estimate continuous performance (recognition) rates of automated speech and gesture recognition systems to inform the models of SIL and GIL , respectively (Figure 4.14).

For speech recognition as a function of human-robot POS , models were trained on annotated spoken phrases from a subset of the data (2 male, 2 female; 25 phrases each; 100 phrase vocabulary). Figure 4.14 illustrates the impact of Equation 3.2 (estimating SOL_{HR} based on POS) on speech recognition rates (SIL_{RH} ; Equation 3.3). The estimation of SOL_{HR} is important for accurately predicting system performance, especially because it often indicates that increase in human vocal effort improves recognition; the alternative is to assume constant speaker behavior, which would predict a performance rate inversely proportional to the distance squared (Wöelfel and McDonough, 2009).

For gesture recognition as a function of human-robot POS , the annotated gesture frames were used to calculate the number of times a particular body part (e.g., head or hands) appeared in the Kinect visual field versus how many times the body part was actually tracked (Figure 4.11). Figure 4.14 illustrates the impact of Equation 3.2 (estimating GOL_{HR} based on POS) on gesture recognition rates (GIL_{RH} ; Equation 3.3) 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.

The 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, alternative techniques were utilized to construct models relating orientation to SIL_{RH} and GIL_{RH} . As noted in Section 3.3.1, two orientations were considered: robot-to-human orientation (α) and human-to-robot orientation (β).

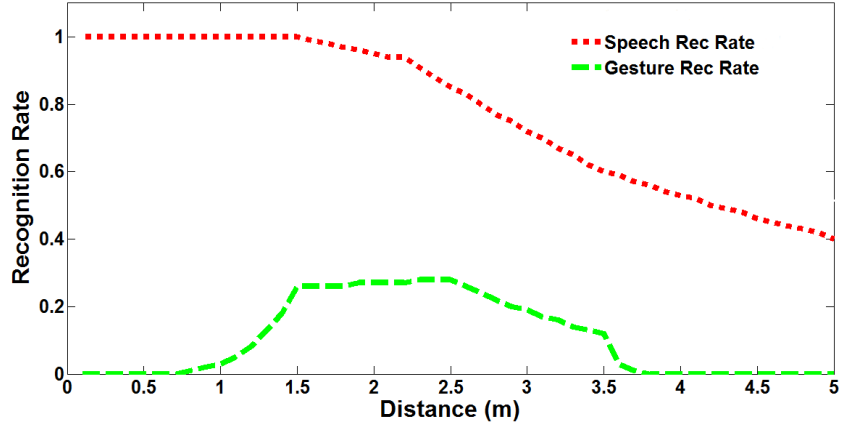


Figure 4.14: Speech and gesture recognition rates as a function of distance.

The impact of robot listener orientation (α) on SIL_{RH} is modeled as the *head-related transfer function* (HRTF) of the PR2 robot, determined using the technique described in Gardner and Martin (1994). The impact of human speaker orientation (β) is based on existing validated models of human *speaker directivity* (Chu and Warnock, 2002). Both models of HRTF and speaker directivity are accessed via lookup tables—the standard approach used in the audio processing community—as there are no closed-form solutions for these functions (Wöelfel and McDonough, 2009). Figure 4.15 illustrates the relationship between orientations (α and β) and SIL_{RH} .

A small, controlled data collection was performed to model the relationship between robot observer orientation (α) and gesture recognition rates. Ten participants (recruited from the Interaction Lab⁷ at the University of Southern California) stood at specified positions and orientations (β), and moved their limbs around in their kinematic workspace. GIL_{RH} was modeled based on human body features tracked by the Kinect (Mead and Matarić, 2012) (Figure 4.11). As with speech input levels (SIL_{RH}), the resulting models are accessed via a lookup tables, as there are no closed-form solutions. Figure 4.16 illustrates the relationships between orientations (α and β) and GIL_{RH} .

⁷<http://robotics.usc.edu/interaction>

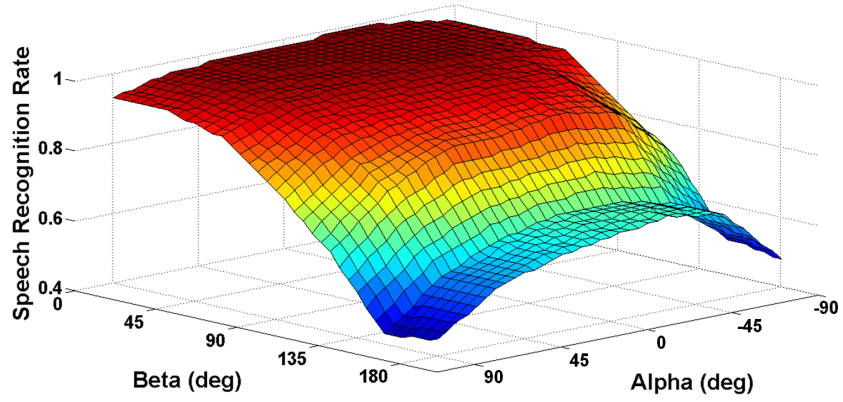


Figure 4.15: Speech recognition rates as a function of human speaker orientation (β) and robot listener orientation (α).

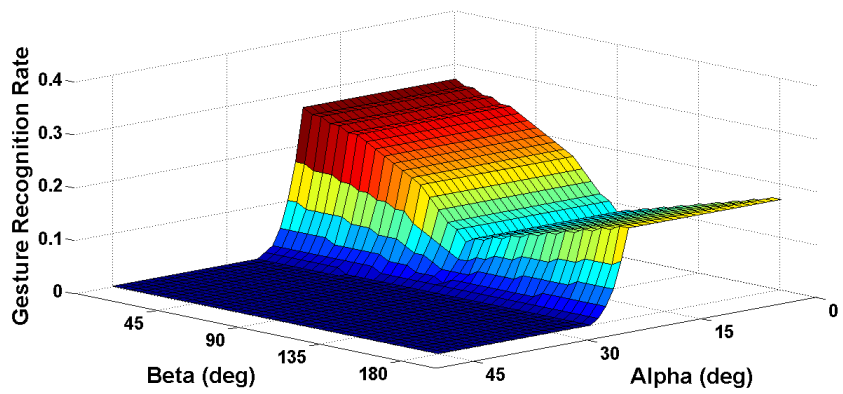


Figure 4.16: Gesture recognition rates as a function of human speaker orientation (β) and robot observer orientation (α).

4.3.3 Extension: Adaptation in Complex Environments

The data collections conducted in this work were performed in a controlled experimental setting; however, real-world scenarios include multiple sources of extrinsic environmental interference, including loud noises or visual occlusions (Adams and Zuckerman, 1991). An autonomous robot might encounter such dynamic and unstructured factors in social interactions. To address this, an extension of the data-driven framework integrates an “**interference**” parameter (INT) into the Bayesian network of proxemics and multi-modal communication (Figure 4.17), extending Equations 3.1 – 3.3 to model the impact of interference (INT) on existing model parameters of: 1) interagent pose (POS); 2) human speech and gesture production (SOL_{HR} and GOL_{HR} , respectively); and 3) robot speech and gesture recognition (SIL_{RH} and GIL_{RH} , respectively). Formally, each component of the model can be written respectively as:

$$p(POS|INT) \tag{4.1}$$

$$p(SOL_{HR}, GOL_{HR}|POS, INT) \tag{4.2}$$

$$p(SIL_{RH}, GIL_{RH}|SOL_{HR}, GOL_{HR}, POS, INT) \tag{4.3}$$

An exhaustive data collection would be necessary to more formally inform the parameters of the models in Equations 4.1 – 4.3. The data collection could employ the same procedures used for the initial modeling of the data-driven approach (described in Section 4.3.1), ensuring compatible integration of the new framework parameters conditioned on INT .

For this dissertation, the parameters of the models are informed by related literature.

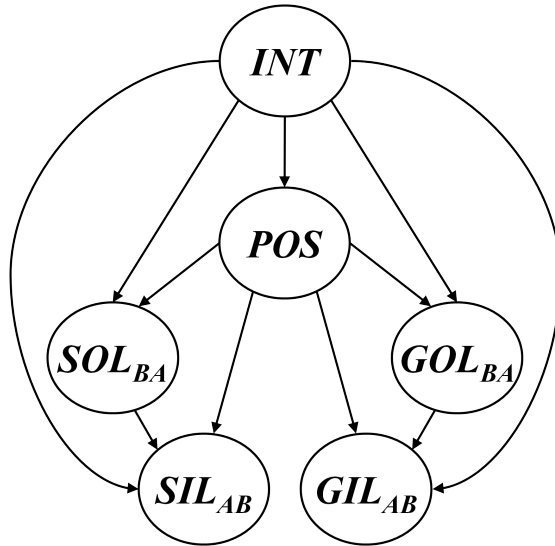


Figure 4.17: A Bayesian network modeling relationships between extrinsic (environmental) interference, pose, speech, and gesture.

Pearsons et al. (1977) extensively investigates the influence of ambient acoustic noise in the environment (INT) on human speech output levels (SOL_{HR}), suggesting that SOL_{HR} increases at a rate of 0.6 db SPL (normalized to 1 meter) per dB SPL increase in background noise (between 48 and 70 dB SPL). Robot speech input levels (SIL_{RH}) remain unchanged, as the models in Section 4.3.2.3 are based on the *relative* (rather than *absolute*) sound pressure level in the environment (i.e., a signal-to-noise ratio). Ozyurek (2002) suggests that humans will adapt their gesture output levels (GIL_{HR}) to accommodate the needs of the observer (in this case, the robot) based on visual occlusions (INT). However, the Kinect-based robot gesture input levels (GIL_{RH}) utilized in this work fails if the head or either shoulder of the human is occluded; thus, for practical purposes, the models of GOL_{RH} predict an inability to recognize gestures (i.e., a value of 0.0) at any interagent pose (POS) in which more than 70% of the body of the human is occluded.

The integration of *INT* as a parameter of the model is a novel consideration in the field of HRI, and illustrates the extensibility and situatedness of the framework. In complex environments, the extended model enables the robot to react to extrinsic interference—such as loud noises (detected via on-board microphones) and visual occlusions (detected via on-board cameras, distance sensors, or an internal map)—and appropriately adapt its proxemic behavior (by moving closer or farther away) and multimodal communication (by speaking louder or gesturing in different regions of space).

4.3.4 Discussion

The fundamental insight of the data-driven approach is in the application of its models of human-robot proxemics to enable situated multimodal communication in HRI. Autonomous sociable robots must utilize automated recognition systems to reliably recognize natural human speech and body gestures (Breazeal, 2003, 2004). The reported models of human multimodal communication (speech and gesture) enable the system to predict the manner in which a social signal will be produced by a person (Section 4.3.2.2, Figure 4.13), which can then be used to predict robot automated social signal recognition rates (Section 4.3.2.3, Figure 4.14). For example, if the robot can detect its current interagent pose (*POS*), then it can use its models to predict 1) how loudly a person will likely speak, 2) how loudly its microphones will likely detect the speech, and 3) how well its automated speech recognition system is likely to perform—all before a single word is spoken by the person. If the sociable robot is mobile, it can 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 (implemented and evaluated in Section 5.3). This is a fundamental capability that autonomous sociable robots should have (Breazeal, 2004), and has the potential to improve autonomy, increase richness of interactions, and generally make robots more reliable, adaptable, and usable in HRI.

4.4 Summary

This chapter discussed the processes of modeling proxemics and multimodal communication using both the heuristic and data-driven approaches presented in this dissertation.

Models of the heuristic approach illustrated that errors in individual and physical feature extraction were a product of the sensor technology used, and that errors in psychological and psychophysical feature extraction were a function of the feature interval size. Furthermore, the encodings of psychophysical features in the heuristic approach are noted to be representative of the human sensory system, which does not accurately extend to robot sensors and recognition systems; thus, while potentially useful for automated human behavior recognition (implemented and evaluated in Section 5.2), the heuristic approach is inadequate for autonomous robot behavior production—this necessitates the models of the data-driven approach.

In the data-driven approach, data were collected and analyzed to inform probabilistic graphical models of proxemics and multimodal communication in human-human and human-robot interactions. The resulting models answer the questions:

1. “Where will a person likely stand with respect to another sociable agent (human or robot)?”
2. “How will a person likely *produce* social signals in different proxemic configurations?”
3. “How will a robot *perceive* those social signals?”

The models of the data-driven approach have implications for the design of robust proxemic behavior and multimodal communication control systems for autonomous sociable robots (implemented and evaluated in Section 5.3).

The next chapter discusses implementations and evaluations of systems based on the heuristic and data-driven models of proxemics and multimodal communication.

Chapter 5

Implementation and Evaluation

This chapter discusses the implementations of both the heuristic and data-driven approaches for unifying proxemics and multimodal communication into deployable recognition and control systems for human-human and human-robot interactions. Three studies are presented to evaluate the performance and user acceptance of the developed systems. Implications of the systems are discussed in the context of autonomous sociable robots.

5.1 Overview

To investigate the efficacy of the computational framework of situated proxemics and multimodal communication established in this dissertation, the approaches discussed— heuristic and data-driven—were integrated into autonomous systems for sociable robots. This chapter presents three studies that implement the approaches as automated human behavior recognition and autonomous robot behavior generation systems, and evaluate these systems with respect to objective measures (system performance), subjective measures (user perceptions), and behavioral measures (changes in user behavior). These three studies are summarized below.

Section 5.2 describes the implementation of the heuristic approach for autonomous human social behavior recognition. Probabilistic graphical models are trained to classify two human interaction cues—*initiation* and *termination*—based on three proxemic feature representations—*physical*, *psychological*, and *psychophysical*. The classification accuracy of each representation (objective measures) is compared.

Section 5.3 describes the implementation of the data-driven approach for autonomous robot proxemic behavior and multimodal communication generation. The system selects pose, speech, and gesture parameters to maximize social signal perception by both the robot and the human user. The system is evaluated with respect to robot automated speech and gesture recognition rates (objective measures).

Section 5.4 evaluates the overall psychophysical philosophy of the computational framework. This study investigates how users perceive robots that vary in performance (subjective measures), and how users adapt their own proxemic behaviors (behavioral measures) to improve automated robot social signal recognition (objective measures).

5.2 Study 1: Human Behavior Recognition

An interaction study was conducted to demonstrate the efficacy of the heuristic approach in the real-time annotation of proxemic features, and to demonstrate its utility in recognizing higher-order human spatiotemporal behaviors—initiation and termination of interaction—in multi-party social encounters.

5.2.1 Data Collection

In the study, two participants—*P1* and *P2*—engaged in an interaction loosely focused on a common object of interest—a static, non-interactive upper-body humanoid robot. The study design, setup, procedures, and resulting dataset are described below.

5.2.1.1 Design

The study procedure (described below in Section 5.2.1.3) was designed to capture natural human proxemic behaviors signifying transitions into (initiation) and out of (termination) social interactions (Deutsch, 1977). An *initiation cue* is a behavior that attempts to engage a potential social partner in discourse, also referred to as a “sociopetal” cue (Lawson, 2001, Low and Lawrence-Ziga, 2003). A *termination cue* proposes the end of an interaction in a socially appropriate manner, also referred to as a “sociofugal” cue (Lawson, 2001, Low and Lawrence-Ziga, 2003, Sommer, 1967). These cues are often targeted (either consciously or subconsciously) at a social partner or stimulus (i.e., an object or another agent, in this case, a robot), and can occur in sequence or in parallel.

5.2.1.2 Setup, Materials, and Measures

The study was set up and conducted in a 6m-by-6m room in the Interaction Lab at the University of Southern California (Figure 5.1). A large physical divider separated the two participants (i.e., they could not see each other) at the beginning of the study; neither participant knew that the other was also participating in the study. The interaction was monitored by the PrimeSensor markerless motion capture system, as well as an overhead color camera and an omnidirectional microphone.

The PrimeSensor structured light range sensor and the OpenNI¹ person tracking software were used for human markerless motion capture. The sensor provides color and depth images, as well as the 3D position and confidence of 24 joints² for each tracked participant. The sensor captured timestamped (for synchronization) images at 320-by-240 images and joint positions at a frequency of 60 Hz.

¹<http://wiki.ros.org/openni>

²The following joints were tracked for each participant: head, neck, torso, waist, collars (left and right), shoulders (left and right), elbows (left and right), wrists (left and right), hands (left and right), longest finger tips (left and right), hips (left and right), knees (left and right), ankles (left and right), and feet (left and right).

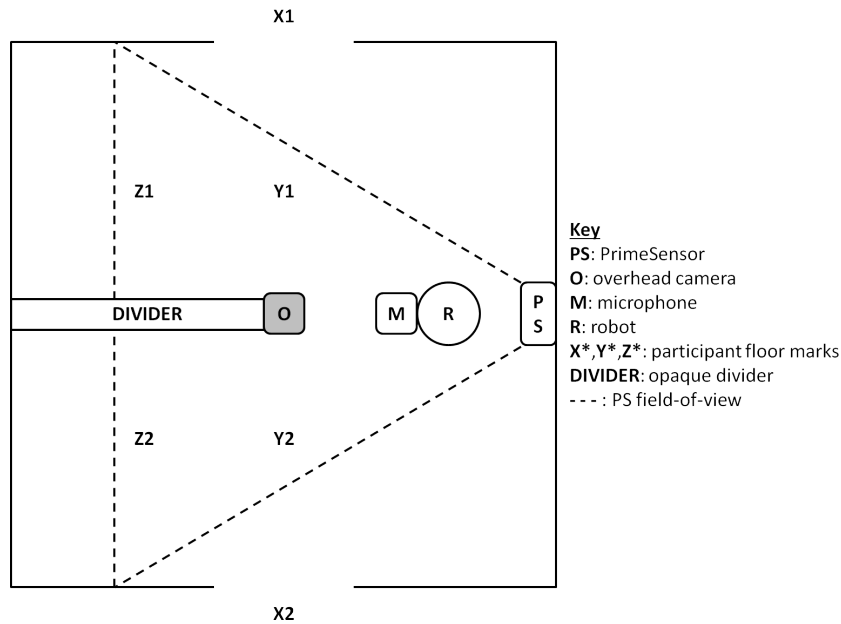


Figure 5.1: The experimental setup for eliciting and recognizing human initiation and termination cues.

The PrimeSensor is capable of tracking people without user-specific calibration; however, this method was determined to be unsuitable for the purposes of this study. Instead, a noninvasive “initialization pose” was used to calibrate the sensor to each person in its field-of-view. Each participant was trained, prior to the execution of the study, to stand straight up, extend his or her arms bilaterally from his or her body, and bend his or her arms at the elbow to form a 90 angle between the upper arm and the forearm. The calibration process took 2–5 seconds.

Each participant was informed that the experimenter would shine a laser pointer at his or her feet to indicate different phases of the experiment: at the beginning, the laser pointer indicated for the participant to approach the robot and begin interacting with it (or anyone else who was there); later, the laser pointer indicated the end of the interaction, and participants would depart the experimental space.

The next section describes the study procedures within this experimental setup.

5.2.1.3 Procedures

P1 silently entered the room first from floor mark *X1*, and stood on floor marks *Y1* then *Z1* for user-specific Primesensor calibration. Shortly thereafter, *P2* silently entered the room from floor mark *X2*, and performed the same calibration at floor marks *Y2* then *Z2*. Note that, from all marked participant locations, the physical divider obstructed each participant's view of the other participant (i.e., participants could neither see nor were they aware that the other participant was in the room; follow-up discussions with each participant confirmed this).

After the calibration, the experimenter shined a laser pointer at the feet of *P1*, indicating that he or she should approach the robot. When *P1* moved away from floor mark *Z1* and approached the robot (an initiation cue directed at the robot), the interaction scenario was considered to have officially begun. Once the participant verbally engaged the robot (unaware that the robot would not respond), *P2* was signaled (via laser pointer behind the field-of-view of *P1*) to approach the dyad from behind the divider, and enter the existing interaction between *P1* and the robot (an initiation cue directed at both *P1* and the robot, often eliciting an additional initiation cue from *P1* directed at *P2*). After joining the group, the participants engaged in an open-ended (i.e., unscripted) conversation that lasted approximately 5–6 minutes. After this period of time, the experimenter shined the laser pointer to indicate that participants should wrap up their conversation and exit the room through the way in which they entered (a termination behavior directed at both the other participant and the robot). Once both participants exited the room, the interaction scenario was considered to be complete.

5.2.1.4 Participants

Participants were recruited via mailing list and word-of-mouth on the campus of the University of Southern California. A total of 36 participants (24 male, 12 female) were involved in the study, paired into 18 interactions. All participants were university students between the ages of 18 and 32, and all had technical backgrounds; 28 of these participants had prior experiences with robots.

5.2.1.5 Dataset

Joint positions recorded by the PrimeSensor were processed in real-time (i.e., at the same 60 Hz frequency as the sensor) to automatically extract individual (Schegloff, 1998), physical (Mehrabian, 1972), psychological (Hall, 1966), and psychophysical (Hall, 1963) proxemic features using the heuristic approach (Sections 3.2 and 4.2).

Images from the PrimeSensor were independently annotated by two coders for occurrences of *initiation* and *termination* cues for each social dyad (pair of agents, *A* and *B*); interrater reliability was high for both initiation ($r_{initiation} = 0.96$) and termination ($r_{gesture} = 0.91$) annotations. The dataset provided 71 usable examples of initiation and 69 usable examples of termination.

This dataset is publicly available upon request.

5.2.2 Data Modeling, Analysis, and Results: Objective Measures

To examine the utility of the heuristic approach for automated human behavior recognition, features extracted from the data collected in the study were used to train Hidden Markov Models (HMMs; (Rabiner, 1990)) to detect interaction cues—*initiation* and *termination*. Section 5.2.2.1 provides an overview of HMMs and how they were used to model the interaction cues, and the resulting performance of the HMMs is presented in Section 5.2.2.2.

5.2.2.1 Hidden Markov Models for Human Behavior Recognition

In this dissertation, five-state left-right HMMs with two skip-states were used (Figure 5.2); this is a common topology used in recognition applications (Rabiner, 1990). Observation vectors for each of the physical, psychological, and psychophysical representations consisted of seven, three, and eight features, respectively (Table 5.1). For each representation, a separate HMM was trained for each of two labels: one for *initiation* cues, and one for *termination* cues. Given observation data, each model returned the likelihood of the data being labeled as the corresponding interaction cue (initiation or termination). The observation and transition parameters were processed and converged after six iterations of the Baum-Welch algorithm (Dempster et al., 1977). Leave-one-out cross-validation was utilized to validate the performance of the models.

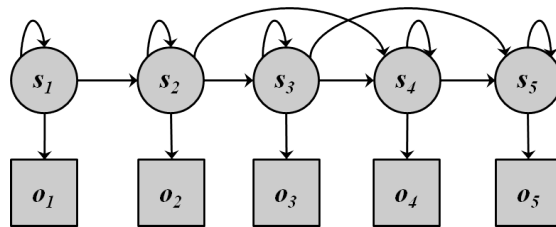


Figure 5.2: A five-state left-right HMM with two skip-states (Rabiner, 1990) used to model each modeled interaction cue—*initiation* or *termination*—for each representation.

Physical Features	Psychological Features	Psychophysical Features
Total Distance	Distance Code	Visual Code (A-to-B)
Straight-Ahead Distance (A-to-B)	SFP Axis Code (A-to-B)	Visual Code (B-to-A)
Straight-Ahead Distance (B-to-A)	SFP Axis Code (B-to-A)	Voice Loudness Code
Lateral Distance (A-to-B)		Kinesthetic Code (A-to-B)
Lateral Distance (B-to-A)		Kinesthetic Code (B-to-A)
Relative Body Orientation (A-to-B)		Olfaction Code
Relative Body Orientation (B-to-A)		Thermal Code
		Touch Code

Table 5.1: The observation vectors for 7-dimensional physical features, 3-dimensional psychological features, and 8-dimensional psychophysical features.

5.2.2.2 Results

Table 5.2 presents the confusion matrix of the HMMs trained on the physical, psychological, and psychophysical feature sets, and Figure 5.3 compares the overall classification accuracy of each representation. The performance of each approach is discussed below.

The HMMs trained using the *physical* features (Mehrabian, 1972) were shown to discriminate between initiation and termination cues (Table 5.2); however, there was often misclassification, resulting in an overall classification accuracy of 56% (Figure 5.3). This is due to the inability of the physical features to capture the complexity of the environment and its impact on the agent’s perception of social stimuli (in this case, the visual obstruction/divider between the two participants).

The HMMs trained using the *psychological* features resulted in a 52% classification accuracy (Figure 5.3), which is little more than chance (50%) at discrimination initiation from termination (Table 5.2). As with the physical representation, the psychological features do not capture the visual occlusion in the scene, resulting in misclassification. Furthermore, the psychological features are discretized with large feature intervals (for both the distance and SFP axis codes), while the physical representation is continuous and performs slightly better; this discretization introduces uncertainty into the classification.

The HMMs trained using the *psychophysical* features showed considerable improvement over the physical and psychological representations (Table 5.2), with an overall classification accuracy of 72% (Figure 5.3). Psychophysical features account for the sensory experience of each agent (specifically, the visual occlusion of the physical barrier), resulting in a more situated and robust representation. While a larger data collection would likely result in an improvement in the recognition rate of each approach, it is anticipated that the relative performance between them would remain unchanged.

		Physical		Psychological		Psychophysical	
		Initiation	Termination	Initiation	Termination	Initiation	Termination
Actual Label	Initiation	41	30	38	33	47	24
	Termination	32	37	34	35	15	54
		Predicted Label					

Table 5.2: Confusion matrix for HMM-based recognition of initiation and termination cues using physical, psychological, and psychophysical feature representations.

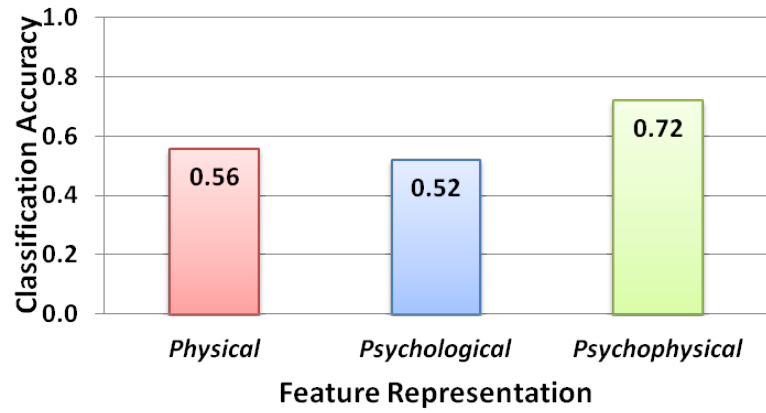


Figure 5.3: Comparison of HMM classification accuracy of initiation and termination behaviors trained over physical, psychological, and psychophysical feature sets.

5.2.3 Discussion

The intuition behind the improvement in performance of the psychophysical representation over the physical and psychological representations is in the incorporation of the sensory experiences of each agent situated within the environment—the psychophysical representation embraces this, whereas the others do not. Specifically, the psychophysical representation encodes the visual occlusion of the physical divider separating participants at the beginning of the interaction. For further intuition, consider two people standing 1 meter apart, but on opposite sides of a door; the physical and psychological representations would mistakenly classify this as an adequate proxemic scenario (based on distance), while the psychophysical representation would correctly classify this as an inadequate proxemic scenario (because the people are visually occluded).

This study demonstrated the feasibility the heuristic approach to autonomously extract physical, psychological, and psychophysical proxemic features in real time, and the utility of the approach for recognizing human spatiotemporal behaviors that signify transitions into (initiation) and out of (termination) a social interaction. The results demonstrated that the HMMs trained on psychophysical features outperform those trained on physical and psychological features suggesting a more powerful representation for recognizing human interaction cues. These results have particular implications for autonomous interactive systems, including (but not limited to) sociable and socially assistive robots. The recognition system is publicly available as open-source software.

While the heuristic approach for proxemics and multimodal communication is useful for automated human behavior recognition, the merits of the psychophysical representation (its ability to predict social stimuli detected in different proxemic configurations) are not accurately modeled for a sociable robot. The next section investigates the implementation and use of the data-driven approach for robot decision-making and control for proxemic and multimodal communication behavior generation.

5.3 Study 2: Robot Behavior Generation

The models (Section 4.3) of the data-driven approach (Section 3.3) of proxemic behavior and multimodal communication were implemented as autonomous systems for sociable robots, described in Section 5.3.1 and analyzed in Section 5.3.2. These systems include:

1. a sampling-based approach for robot *goal state estimation* with respect to proxemic and multimodal communication parameters (presented in Section 5.3.1.1),
2. a *reactive proxemic controller* for goal state realization (presented in Section 5.3.1.2), and
3. a *novel cost-based trajectory planner* for maximizing autonomous robot speech and gesture recognition along a path to the goal state (presented in Section 5.3.1.3).

5.3.1 Proxemic Behavior and Multimodal Communication Systems

The proxemic behavior and multimodal communication systems implemented using the data-driven approach (Section 3.3) and models (Section 4.3) are presented below. Each of these systems is publicly available as open-source software.

Section 5.3.1.1 discusses the decision-making process by which a sociable robot selects a goal state (a set of interagent pose, and speech and gesture output and input parameters) that is optimal for it to engage in social interactions based on the models of predicted sensory experience of all interacting agents (human and robot). The resulting proxemic goal state often differs from that reported in previous work in human-robot proxemics (e.g., to track a person using its on-board markerless motion capture system, the robot typically needs to be farther away than people traditionally prefer), but robot performance (as measured by automated speech and gesture recognition rates) is *objectively* improved (reported in Section 5.3.2 below); human *subjective* and *behavioral* responses to this atypical robot behavior are investigated in Study 3 (Section 5.4).

Section 5.3.1.2 presents a reactive proxemic controller that moves the robot to the desired goal pose. The controller is simple to implement and enables the robot to quickly reach its goal pose; however, the controller can fail if a movement places the robot in location from which it cannot see the human. This failure condition motivated the development of a trajectory planner that encapsulates this constraint.

Section 5.3.1.3 addresses potential shortcomings of the reactive controller by developing a novel cost-based trajectory planner that considers “interaction potential”, a measure of how well the robot would recognize multimodal communication signals produced by the human while the robot is moving to its proxemic goal state.

Section 5.3.2 provides an analysis of the implemented systems with respect to *objective* measures of predicted robot speech and gesture recognition rates, as well as path lengths traversed by the robot.

5.3.1.1 Goal State Estimation for Pose, Speech, and Gesture Parameters

The data-driven models of human communication reported in 4.3.2.2 enable the system to predict the manner in which multimodal communication will be produced by a person (Figure 4.13), which can then be used to predict robot performance as measured by automated speech and gesture recognition systems (Section 4.3.2.3, Figure 4.14). Given an interagent pose (POS), the robot can use these models to predict 1) how the user will likely speak and gesture (SOL_{HR} and GOL_{HR} , respectively), 2) how its sensors will likely detect the speech and gestures, and 3) how well its speech and gesture recognition systems will likely perform (SIL_{RH} and GIL_{RH} , respectively)—all before the person starts speaking or gesturing.

These predictions inform the development of a decision-making mechanism for a robot to select (and subsequently move to) a better pose (POS) to maximize the potential for its performance in the interaction; furthermore, the same mechanism can be used to select robot speech and gesture output parameters (SOL_{RH} and GOL_{RH} , respectively) to maximize human understanding of its own social signals. Formally, this decision-making mechanism is used for robot *goal state estimation*, where the robot goal state is defined by three parameters:

1. **robot interagent pose**, POS —Where should the robot R position itself relative to the human H ?
2. **robot speech output level**, SOL_{RH} —How loudly should the robot R speak for the human H to hear?
3. **robot gesture output level**, GOL_{RH} —In what region of space should the robot R gesture for the human H to see?

Using the Bayesian models developed for the data-driven approach (Section 4.3.2), the goal state estimation system infers information about unobserved variables (e.g., how a person will speak or where a person will gesture) based on sensed data (e.g., current interagent pose), and uses these inferences to select parameters for robot social behavior (e.g., interagent pose, and speech and gesture output levels).

To estimate optimal robot interaction parameters, the system uses a sampling-based approach to consider interagent poses that the robot could occupy (Figure 5.4). At each sampled pose ($POS^{[i]}$), the system predicts how all interacting agents would produce (output) and perceive (input) social stimuli (speech and gesture) if the interaction occurred with the robot at that pose; robot speech and gesture output levels corresponding to $POS^{[i]}$ ($SOL_{RH}^{[i]}$ and $GOL_{RH}^{[i]}$, respectively) are estimated based on the models of human-human multimodal communication reported in Section 4.3.2.2 (i.e., the robot

modulates its speech and gestures in the same way that a human does). The pose is then assigned a weight ($w^{[i]}$), evaluating the likelihood of social signal perception by all interacting agents. The interaction pose and robot speech/gesture output levels with the maximum weight are selected as the goal state parameters (Equations 5.1 and 5.2, respectively). Conditioning these parameters on the previous state during resampling ensures that they will not change drastically between time intervals.

$$\operatorname{argmax}_{POS} E[SIL_{RH}, GIL_{RH} | SOL_{HR}, GOL_{HR}, POS] \quad (5.1)$$

$$\operatorname{argmax}_{SOL_{RH}, GOL_{RH}} E[SIL_{HR}, GIL_{HR} | SOL_{RH}, GOL_{RH}, POS] \quad (5.2)$$

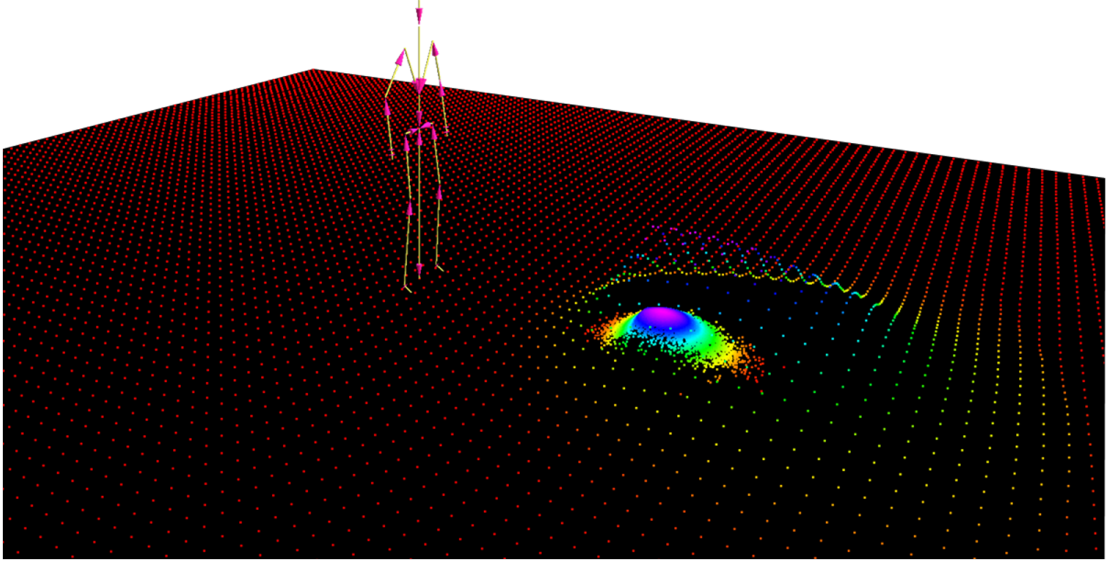


Figure 5.4: The goal state estimation system uses a sampling-based approach to estimate how the human produces social signals, and how the robot perceives them using models in the data-driven computational framework of proxemics and multimodal communication. Low to high estimates are denoted by both the coloring (red to magenta, respectively) and height (low to high, respectively) of sampled points in the floor plane (shown in black). The clustering of points denotes the impact of the resampling process. The system uses these estimates to select parameters that maximize the expected performance of the robot during the interaction.

Once appropriate goal state parameters— POS , SOL_{RH} , and GOL_{RH} —have been selected, the robot must execute actions to reach the goal state. For SOL_{RH} and GOL_{RH} , this is done by adjusting the robot speaker volume and gesture locations, respectively (Mead et al., 2010). For POS , this is achieved using either a reactive controller or a trajectory planner, described in Sections 5.3.1.2 and 5.3.1.3, respectively.

5.3.1.2 Reactive Proxemic Controller

To reach the desired pose provided by the goal state estimation system (Section 5.3.1.1), a reactive proxemic controller was developed for autonomous sociable robots.

The reactive proxemic controller accepts the following input signals: desired and actual interagent distances (D and d , respectively); desired and actual angle from the robot to the person (A and α , respectively), and the desired and actual angle from the person to the robot (B and β , respectively). The desired distance and angles (D , A , and B) come from the proxemic goal state, and the actual distance and angles (d , α , and β) are extracted from sensor readings. The controller attempts to minimize the error between each of the desired and actual distance and angle signals; the resulting linear velocities (v_x , v_y) and angular velocity (v_ω) are expressed as follows:

$$\begin{pmatrix} v_x \\ v_y \\ v_\omega \end{pmatrix} = \begin{pmatrix} \cos(\alpha) & -\sin(\alpha) & 0 \\ \sin(\alpha) & \cos(\alpha) & 0 \\ 0 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} (d - D) \times \cos(\alpha - A) \\ \sin(\beta - B) \\ \alpha - A \end{pmatrix} \quad (5.3)$$

The reactive proxemic controller was implemented on the PR2 robot platform from Willow Garage. The PR2 features a semi-holonomic mobile base, which affords it the ability to strafe (v_y), a valuable motion for proxemic control of the “lateral distance” feature in the physical representation (Section 3.2.2). Interagent pose data (d , α , and β) were extracted using an on-board Microsoft Kinect (described in Section 4.3).

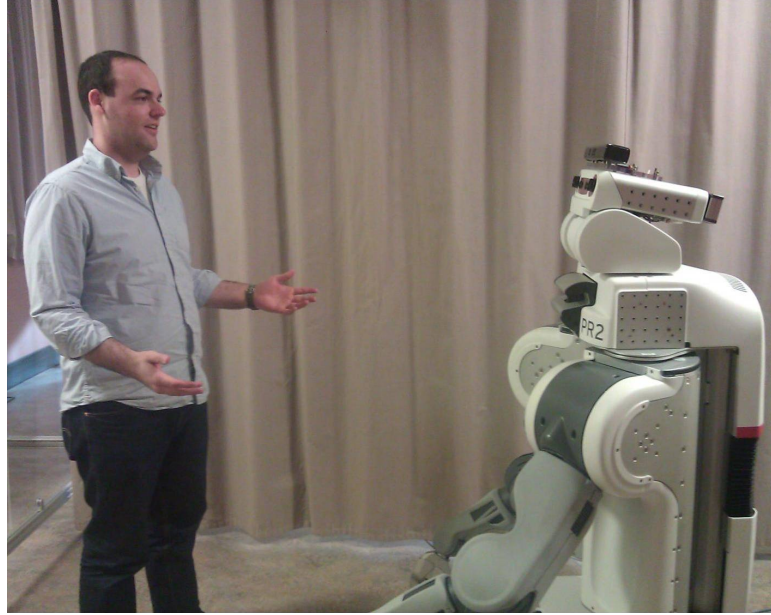


Figure 5.5: The PR2 robot utilizes the reactive proxemic controller to reach the desired pose provided by the goal state estimation system during an interaction with a user.

A preliminary evaluation of the controller was performed with a small group of human users. Initial tests revealed that the controller successfully reached the desired proxemic pose when the robot started its approach from a location in front of the user (i.e., $|\beta| < 90^\circ$, a region in which interactions typically occur), but sometimes failed when the robot started its approach from a location behind the user (i.e., $|\beta| > 90^\circ$). This occurred because the controller does not consider the ability of the robot to track the user (part of the model prediction of robot gesture input level, GIL_{RH} ; Section 4.3.2.3) during its motion—thus, at any time during its movement to the goal pose, if the robot reaches a location from which it can no longer track the user, the robot halts its motion, as the goal pose (used to calculate robot velocity) is relative to and necessitates a tracked human user (a complete analysis is provided below in Section 5.3.2.2). This insight led to the development of a trajectory planner that considers the sensory experience of the robot along its path to the goal pose, presented in the next section.

5.3.1.3 Cost-based Trajectory Planner using Interaction Potential

To alleviate the shortcomings of the reactive proxemic controller, a cost-based trajectory planner was implemented to reach the goal pose provided by the goal state estimation system (Section 5.3.1.1). Traditional cost-based planners typically weight candidate trajectories based on a set of distance-based heuristics, such as distance to the goal pose and distance to nearby obstacles (Marder-Eppstein et al., 2010); however, these planners do not consider the ability of the robot to track a human user (who is potentially moving) at each location along the path. The dissertation addresses this by further weighting each candidate trajectory based on a metric of “interaction potential”, representing a prediction of how the robot would perceive social stimuli (speech and gesture) produced by a user if the interaction were to occur at any moment in time along its path.

Trautman and Krause (2010) introduced the term “interaction potential” in the context of robot navigation in dense crowds of people, defining it as a likelihood of a mobile robot entering the personal space of a human (who might also be moving) while the robot is trying to reach a goal location within a shared environment. In this related work, interaction potential was evaluated based on the distance to each person (i.e., a single *physical* feature, as defined in Section 3.2), and robot trajectories were selected to *minimize* interaction potential (i.e., the robot safely avoided people).

This dissertation considers interaction potential from a multi-feature *psychophysical* perspective (Section 3.3), evaluating the predicted robot speech and gesture recognition rates (SIL_{RH}^t and GIL_{RH}^t , respectively) at each pose (POS^t) at a time t along a candidate path (based on Equation 3.3), with the goal of *maximizing* interaction potential along the path to the goal pose. The interaction potential (IP) for a given candidate velocity (v_x^t, v_y^t, v_ω^t) at time t along a projected path is defined as:

$$IP(v_x^t, v_y^t, v_\omega^t) = p(SIL_{RH}^{t+1}, GIL_{RH}^{t+1} | SOL_{HR}^{t+1}, GOL_{HR}^{t+1}, POS^{t+1}) \quad (5.4)$$

Related work in sociable robot navigation has demonstrated the feasibility of real-time person-aware trajectory planning by weighting a cost function based on its fit to a more socially appropriate model of navigation (Feil-Seifer and Matarić, 2011b). A similar approach is employed by this dissertation. Equation 5.4 is used to extend an existing cost function ($cost$), establishing a new cost function ($cost_{IP}$) weighted by IP^3 (w_{IP}^t):

$$w_{IP}^t = IP^{-1}(v_x^t, v_y^t, v_\omega^t) \quad (5.5)$$

$$cost_{IP}(v_x^t, v_y^t, v_\omega^t) = w_{IP}^t \times cost(v_x^t, v_y^t, v_\omega^t) \quad (5.6)$$

Thus, for $IP \approx 1$, the resulting IP -weighted cost function remains nearly identical to the unweighted cost function (i.e., $cost_{IP} \approx cost$). Trajectories in which the robot cannot interact with the user are penalized⁴, as shown by $\lim_{IP \rightarrow 0} cost_{IP} = \infty$. This results in safe trajectories that support the potential for the robot to autonomously recognize speech and gestures while moving to its goal pose (Figure 5.6).

The IP -weighted cost-based trajectory planner was implemented by extending the *navigation*⁴ software stack (Marder-Eppstein et al., 2010) within the Robot Operating System (ROS⁵) framework (Quigley et al., 2009). The software uses a *global* planner for planning the entire path from the current pose to the goal pose, and a *local* planner for dynamic replanning around obstacles along the way (Fox et al., 1997, Gerkey and Konolige, 2008). For this dissertation, the cost function employed by the both the global and local planners were weighted as per Equation 5.6; the local planner was also modified to be executed if the person (and, thus, the relative goal pose) moved. As with the reactive proxemic controller (Section 5.3.1.2), the PR2 robot was selected as the hardware platform for testing of the IP -weighted cost-based trajectory planner.

³In practice, a small number (in this dissertation, 10^{-6}) is added to IP to prevent division-by-zero errors when calculating the weight w_{IP}^t (Equation 5.5).

⁴<http://wiki.ros.org/navigation>

⁵<http://www.ros.org>

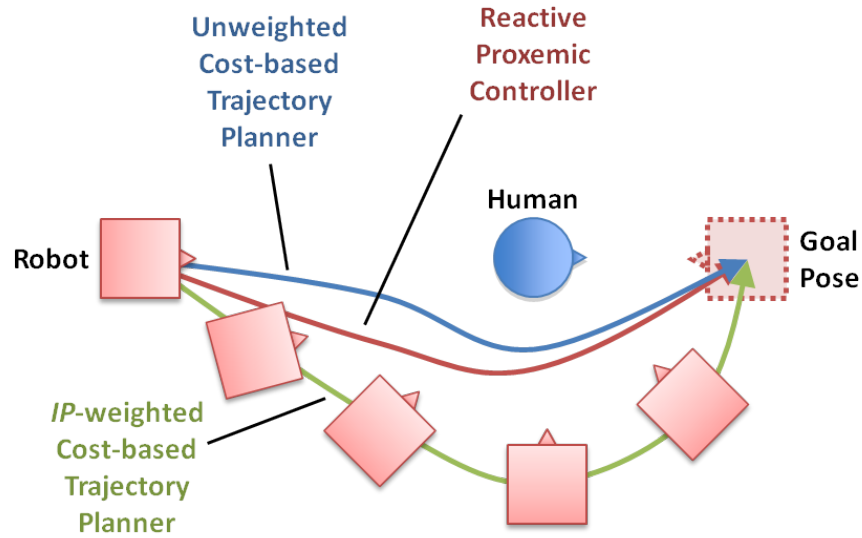


Figure 5.6: Trajectories of an autonomous robot to a goal pose using three different motion control systems: 1) the reactive proxemic controller (red line; Section 5.3.1.2); 2) the *IP*-weighted cost-based trajectory planner (green line; Section 5.3.1.3); and 3) an unweighted cost-based trajectory planner (blue line) (Marder-Eppstein et al., 2010).

5.3.2 System Analysis: Objective Measures

This section provides an analysis of the implemented robot behavior generation systems with respect to *objective* measures of autonomous robot performance based on automated speech and gesture recognition rates. In Section 5.3.2.1, the proxemic configuration predicted by the goal state estimation system (based on robot performance) is compared to models of human-robot proxemics that rely on human preferences. Section 5.3.2.2 compares the reactive proxemic controller (Section 5.3.1.2), the cost-based trajectory planner weighted by interaction potential (Section 5.3.1.3), and an unweighted cost-based trajectory planner with respect to path length, as well as the average interaction potential as the robot traverses the path to the goal pose.

5.3.2.1 Analysis of Goal State Estimation System

The desired robot goal pose provided by the goal state estimation system is compared to the human-preferred pose determined in Section 4.3.2.1. In an dyadic (i.e., two agent—human and robot) interaction under typical environmental conditions (i.e., very little noise and no visual occlusions), the desired pose distance D from the goal state estimation system (considering robot performance) and human preferences (with respect to the PR2 robot) were 1.64 meters and 0.94 meters, respectively; the desired robot-to-human and human-to-robot orientations (A and B , respectively) were both 0° in each model, so it is not necessary to compare them. These distances from the two models—*robot performance* and *human preference*—are evaluated with respect to:

1. **Human-Preferred Interagent Pose:** $p(POS_d = D)$ (Section 4.3.2.1),
2. **Robot Speech Input Level:** $E[SIL_{HR}|POS_d = D]$ (Section 4.3.2.3),
3. **Robot Gesture Input Level:** $E[GIL_{HR}|POS_d = D]$ (Section 4.3.2.3), and
4. **Combined Evaluation:** $E[SIL_{HR}, GIL_{HR}|POS_d = D] \times p(POS_d = D)$.

These results of these evaluations are presented in Table 5.3 and discussed below.

	Desired Distance, D (in meters)	Human-Preferred Interagent Pose	Robot Speech Input Level	Robot Gesture Input Level	Combined Evaluation
<i>Robot Performance</i>	1.64	0.50	1.00	0.35	0.18
<i>Human Preference</i>	0.94	0.65	1.00	0.16	0.11

Table 5.3: A comparison of the evaluation of distances predicted by two models—*robot performance* (using the goal state estimation system; Section 5.3.1.1) and *human preference* (reported in Section 4.3.2.1).

The evaluation of distances predicted by models of robot performance (using the goal state estimation system) vs. human preferences (using results from Section 4.3.2.1) provide insights into the overall performance of the data-driven approach of human-robot proxemics (inspired by psychophysical features, such as speech and gesture perception) and traditional approaches for human-robot proxemic behavior (based on physical features, such as distance and orientation). The distance predicted based on robot performance (1.64 meters) was farther away than human preferences (0.94 meters); however, the difference in these distances is within one standard deviation of human preferences ($SD = 0.61$), so users might accept this proxemic configuration (further discussed in Section 5.3.3 and investigated in Section 5.4). Robot speech recognition rates were equal for both models; however, the robot gesture recognition rate at a distance based on models of robot performance was more than double ($2.2\times$) that at a distance based on human preferences, illustrating the potential for a dramatic increase in robot perception and autonomy during face-to-face HRI. Furthermore, the combined evaluation—considering both robot speech and gesture recognition rates, as well as human preferences—highlights a 64% increase in overall performance using models that incorporate robot performance rather than human preferences alone (Table 5.3).

5.3.2.2 Analysis of Reactive Controller and Trajectory Planners

To evaluate the ability of the robot to autonomously reach the goal pose provided by the goal state estimation system, three motion control systems were compared: 1) the reactive proxemic controller (Section 5.3.1.2), 2) a cost-based trajectory planner weighted by interaction potential (Section 5.3.1.3), and 3) an unweighted cost-based trajectory planner (Marder-Eppstein et al., 2010). These systems were evaluated with respect to two *objective* measures: path length and “interaction potential” (defined in Section 5.3.1.3).

A small data collection was performed in which the PR2 robot autonomously navigated to a goal pose predicted by the goal state estimation system; as reported in Section 5.3.2.1, this goal pose (based on robot performance) was $POS = \{1.64\text{m}, 0^\circ, 0^\circ\}$. The robot began its motion from one of twelve starting poses (i.e., there were twelve trials for each motion control system). These starting poses were distributed radially around an open 6m-by-6m room in 30° intervals around and at a distance of 2.5 meters from the center of the room, with the robot facing the center (Figure 5.7). The center represented the location of a human user and was hard-coded as such on the robot’s map of the room, enabling the robot to continuously move relative to this location regardless of whether or not it was actually able to track the person (for evaluation purposes); the experimenter stood at this location for the evaluation, as the *objective* metrics used (described below) did not necessitate a diverse set of human participants.

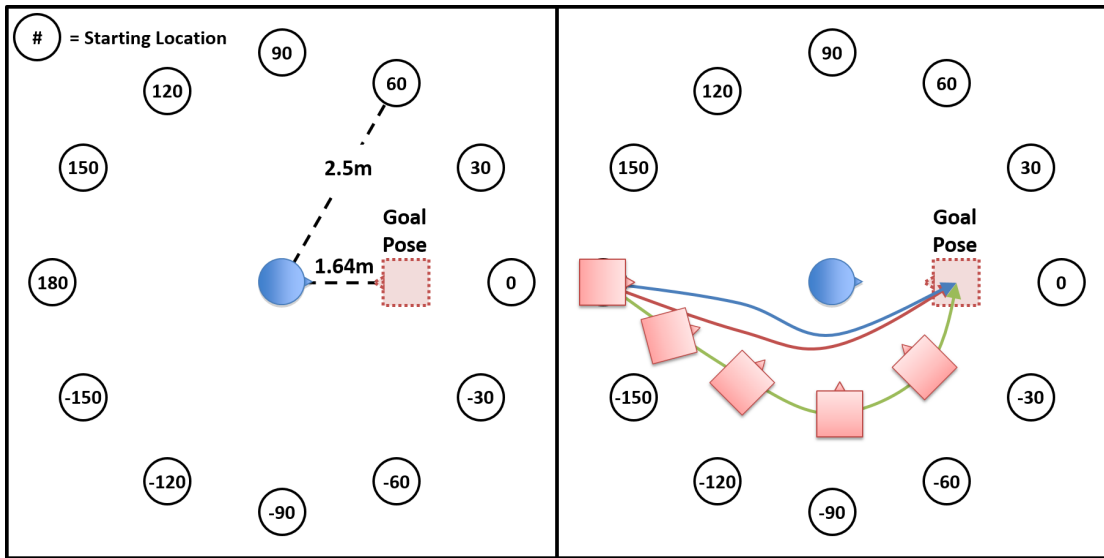


Figure 5.7: The experimental setup (left) and trajectories from the three motion control systems (right): 1) the reactive proxemic controller (red line), 2) the IP -weighted cost-based trajectory planner (green line), and 3) the unweighted cost-based trajectory planner (blue line).

The predicted, simulated (using a simulated PR2 robot in the Gazebo simulator⁶), and measured (using the physical PR2 robot) performance of each motion control system was evaluated with respect to two *objective* measures: 1) *average normalized path length*, and 2) *average interaction potential*. The “average normalized path length” for each motion control system was calculated in three steps: first, the actual path length was determined by the accumulated distance traveled by the robot (as predicted, simulated, or measured by its localization system) as it moved from a starting pose to the goal pose; next, this path length was normalized based on the Euclidean distance from the goal pose to the starting pose corresponding to the trial; finally, the resulting normalized path lengths were averaged for comparison. The predicted and simulated “average interaction potentials” for each motion control system were calculated across trials by averaging the interaction potentials (using Equation 5.4) evaluated at each time t along the path from the starting pose (at $t = 0$) to the goal pose (at $t = t_f$); the measured average interaction potential was determined by counting the number of times that head and hands of the experimenter were detected by the Microsoft Kinect as the robot navigated to the goal pose. The results for each of the motion control systems evaluated with respect to these metrics are presented in Figure 5.8 and discussed below.

⁶<http://gazebosim.org>

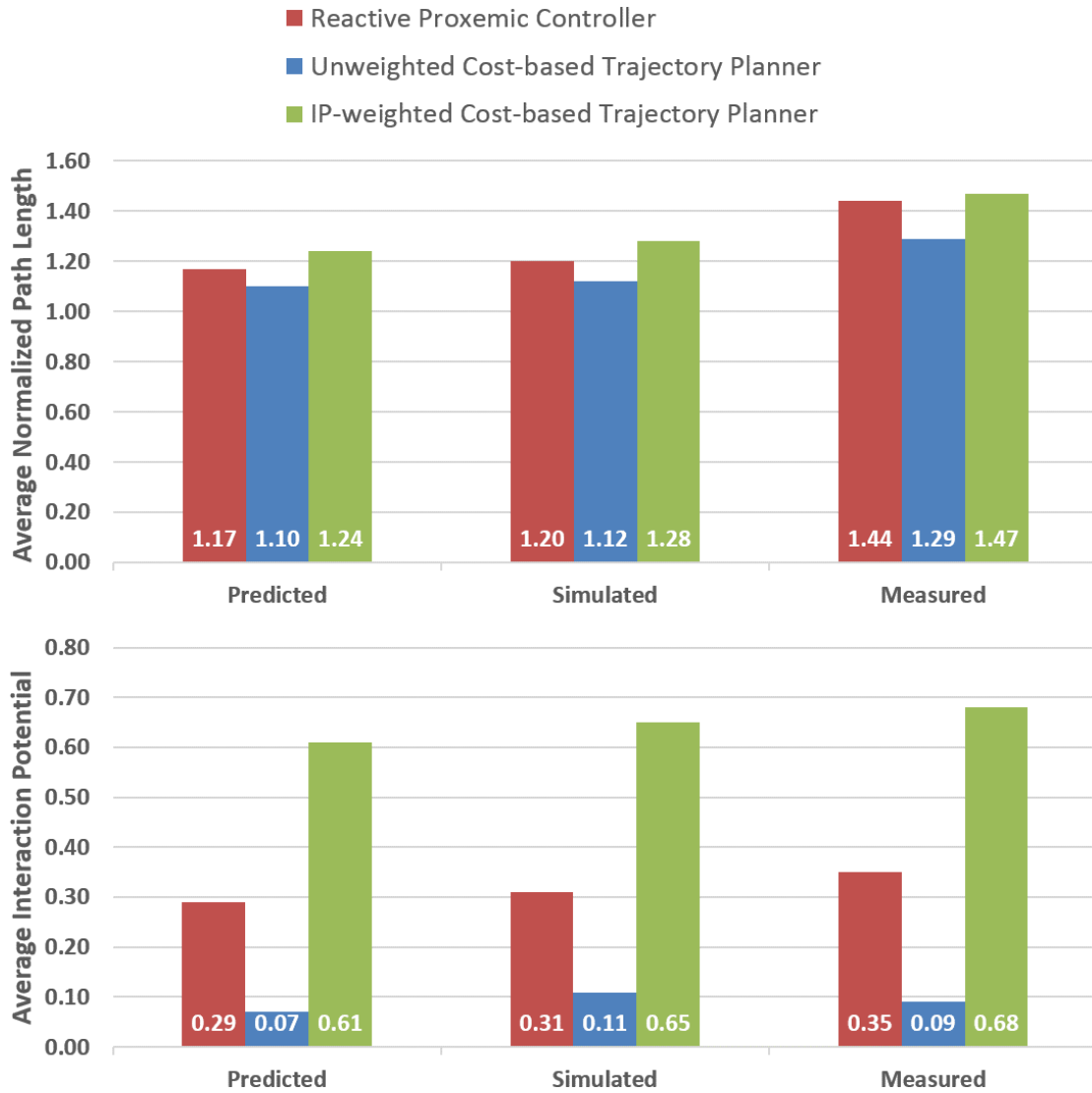


Figure 5.8: A comparison of the *average normalized path length* (top) and *average interaction potential* (bottom) of each of the three motion control systems: 1) the reactive proxemic controller, 2) the *IP*-weighted cost-based trajectory planner, and 3) the unweighted cost-based trajectory planner.

Differences in predicted, simulated, and measured path lengths illustrated in Figure 5.8 were a result of the simulated and actual dynamics of the robot (i.e., it does not perfectly traverse the planned trajectory). These differences, in turn, impacted the interaction potential for each motion control system; measured results were shown to be consistent with those predicted and simulated using Equation 5.4, demonstrating its utility for autonomous sociable robots.

The results highlight that an unweighted cost-based trajectory planner selects paths that minimize the path length to a goal pose (which is its primary function); however, in doing so, the planner will be unable to track a human user (because the robot is either too close or not oriented properly for human recognition), performing at a measured 9% on average, which is unacceptable for autonomous human-robot interactions.

The reactive proxemic controller traversed paths that were marginally longer than that of the unweighted trajectory planner, and performed better on average at detecting the human user (measured at 35%); this increase in performance is a result of error correction in actual vs. desired robot-to-human orientation ($\alpha - A$; Section 5.3.1.2; Equation 5.3), which naturally orients the robot toward the user during its motion, thus, maintaining the user in the Kinect field-of-view for tracking.

The *IP*-weighted cost-based trajectory planner selected paths that were 12–14% longer than the unweighted trajectory planner; however, the utility of the system is demonstrated in its recognition of the human head and hands (measured average of 68%) during robot navigation. Recognition failures of the *IP*-weighted planner—as well as the other motion control systems—occurred when the robot navigated from locations behind ($|\beta| > 90^\circ$) or to the side of ($|\beta| \approx 90^\circ$) the person, as the Kinect person recognition system does not perform well in these configurations (Figure 4.16). Using interaction potential (Equation 5.4), these conditions were predicted, enabling the *IP*-weighted planner to select and traverse the best possible path for interaction.

Experimenter observations from this data collection revealed three general conditions—based on interagent pose parameters (d , α , and β)—in which each of these motion control systems would perform poorly for HRI:

1. If the robot moves to a distance (d) that is too far or too near for speech or gesture recognition, then the system will be unable to interact with a human user.
2. If at any time during navigation the person is located outside of the robot field-of-view (robot-to-human orientation $|\alpha| > 30^\circ$; Figure 4.16), then the robot can neither track the person nor estimate a goal pose with respect to the person, so robot motion ceases movement.
3. If the robot is viewing the person from behind or from the side (i.e., human-to-robot orientation $|\beta| > 90^\circ$ or $|\beta| \approx 90^\circ$, respectively), then automated speech and gesture recognition systems will not perform well (Figures 4.15 and 4.16, respectively).

5.3.3 Discussion

This section presented the application of the data-driven approach (Section 3.3) of proxemic behavior and multimodal communication in autonomous systems for sociable robot behavior generation. A reactive proxemic controller (Section 5.3.1.2) and cost-based trajectory planner based on “interaction potential” (Section 5.3.1.3) were implemented to navigate a robot to an appropriate proxemic pose to conduct interaction. This pose was determined using a goal state estimation system (Section 5.3.1.1) to maximize how the robot would perceive human speech and gestures. The resulting proxemic configuration differed from those reported in related work in human-robot proxemics (Section 2.3.2), which might influence user acceptance of the system. This is investigated in the next section, evaluating the overall psychophysical perspective taken by this dissertation.

5.4 Study 3: Human Acceptance of Robot Behaviors

The results from Study 2 (Section 5.3) demonstrated that the proxemic configuration predicted by the goal state estimation system (Section 5.3.1.1) improved the performance of the robot in recognizing human speech and gestures (based on *objective* measures); however, the controller also produces proxemic configurations that are atypical for human-robot interactions (e.g., positioning itself farther or nearer to the user than preferred). Thus, the question arose as to whether or not people would adopt a technology that places performance over preference, as it might place a burden on people to change their own behaviors to make the technology function adequately.

A follow-up study was performed to investigate how user proxemic preferences change (based on *behavioral* measures) in the presence of a sociable robot that is recognizing and responding to instructions provided by a human user, and how people evaluate (based on *subjective* measures) the robot as its performance varies in different proxemic configurations. Robot performance (ability to understand speech and gesture) is artificially attenuated to expose participants to success and failure scenarios while interacting with the robot, so as to not limit the implications of the results to that of the models developed in this dissertation.

This study represents the largest and most comprehensive experiment conducted in this dissertation. The sections below present the experimental setup, procedure, measures, conditions, and hypotheses, and provide detailed analyses and discussions of behavioral and subjective results.

5.4.1 Experimental Setup

5.4.1.1 Materials

The experimental robotic system used in this work was the Bandit upper-body humanoid robot⁷ (Figure 5.9). Bandit has 19 degrees of freedom: 7 in each arm (shoulder forward-and-backward, shoulder in-and-out, elbow tilt, elbow twist, wrist twist, wrist tilt, grabber open-and-close; left and right arms), 2 in the head (pan and tilt), 2 in the lips (upper and lower), and 1 in the eyebrows. These degrees of freedom allow Bandit to be expressive using individual and combined motions of the head, face, and arms. Mounted atop a Pioneer 3-AT mobile base⁸, the entire robot system is 1.3 meters tall.



Figure 5.9: The Bandit upper-body humanoid robot platform.

⁷<http://robotics.usc.edu/interaction/?l=Laboratory:Robots#BanditII>

⁸<http://www.mobilerobots.com/ResearchRobots/P3AT.aspx>

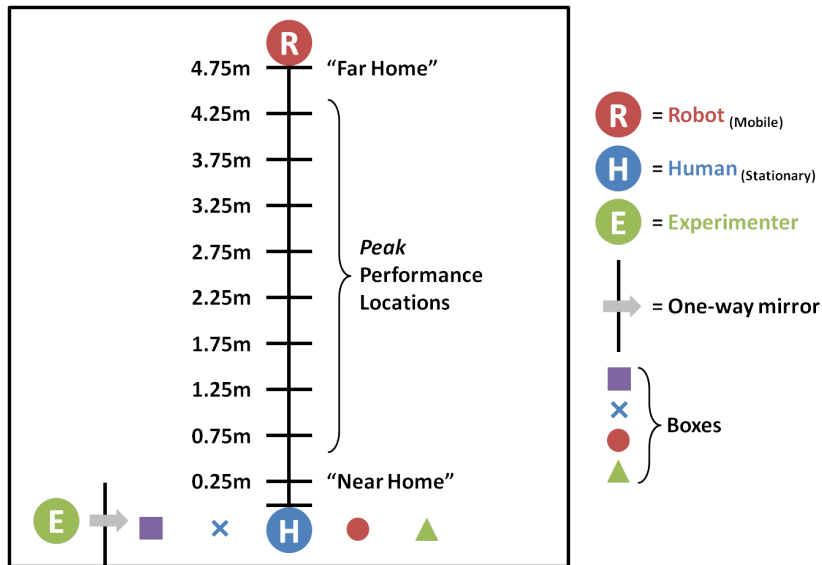


Figure 5.10: The experimental setup for evaluating human acceptance of robot behaviors.

A Bluetooth PlayStation 3 (PS3) controller served as a remote control interface with the robot. The controller was used by the experimenter (seated behind a one-way mirror; Figure 5.10) to step the robot through each part of the experimental procedure (described in Section 5.4.2)—the decisions and actions taken by the robot during the experiment were completely autonomous, but the timing of its actions were controlled by the press of a “next” button. The controller was also used to record distance measurements during the experiment, and to provide ground-truth information to the robot as to what the participant was communicating (however, the robot autonomously determined how to respond based on the experimental conditions described in Section 5.4.3).

Four small boxes were placed in the room, located at 0.75 meters and 1.5 meters from the centerline on each side (left and right) of the participant (Figure 5.10). During the experiment (described in Section 5.4.2), the participant instructed the robot to look at these boxes. Each box was labeled with a unique shape and color; in this experiment, the shapes and colors matched the buttons on the PS3 controller: a green triangle, a

red circle, a blue cross, and a purple square. This allowed the experimenter to easily indicate to the robot to which box the user was attending (i.e., “ground-truth”).

A laser rangefinder on-board the robot was used to measure the distance from the robot to the participant’s legs at all times.

The entire experimental setup is illustrated in Figure 5.10; the experimental procedures within this setup is described in Section 5.4.2. The next section describes the autonomous behaviors of the robot during this data collection.

5.4.1.2 Robot Behaviors

The robot autonomously executed three behaviors throughout the experiment: 1) forward and backward base movement, 2) maintaining eye contact with the participant, and 3) responding to participant instructions with head movements and audio cues.

Robot base movement was along a straight-line path directly in front of the participant, and was limited to distances of 0.25 meters (referred to as the “near home” location) and 4.75 meters (referred to as the “far home” location); it returned repeatedly to these “home” locations throughout the experiment. Robot velocity was proportional to the distance to the goal location; the maximum robot speed was 0.3 m/s, which people find acceptable (Satake et al., 2009).

As the robot moved, it maintained eye contact with the participant. The robot has eyes, but they are not actuated, so the head of the robot was pitched up or down depending on the location of the participant’s head, which was determined by the distance to the participant (from the on-board laser) and the participant’s self-reported height. Prolonged eye contact from the robot has been shown to increase user preferences of distance in HRI (Mumm and Mutlu, 2011, Takayama and Pantofaru, 2009); this was not a problem in the experiment, as the study is concerned with evaluating change in user proxemic preferences rather than absolute values of preference.

The robot provided head movement and audio cues to indicate whether or not it understood instructions provided by the participant (described in Section 5.4.2.3). If the robot understood the instructions, it provided an *affirmative response* (looking at a box); if the robot did not understand the instructions, it provided a *negative response* (shaking its head). With each head movement, one of two affective sounds were also played to supplement the robot’s response; affective sounds were used because robot speech influences proxemic preferences and would have introduced a confound in the experiment (Walters et al., 2008).

Within the described setup, an experiment was performed to investigate user perceptions of robot performance attenuated by distance and its effect on proxemic preferences.

5.4.2 Experimental Procedure and Measures

The experimental procedure was performed in six sequential phases that included a human-robot interaction, and behavioral and subjective measures. Measures are presented inline with the procedure to provide an order of events as they occurred. The sequence of phases is illustrated in Figure 5.11; each phase is discussed below.

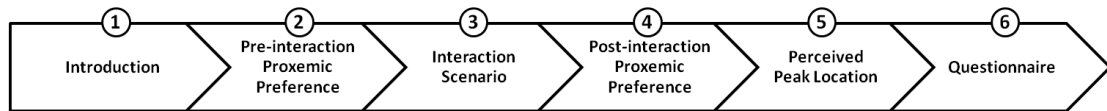


Figure 5.11: The six phases of the experimental procedure.

5.4.2.1 Phase 1: Introduction

Participants (described in Section 5.4.5) were greeted at the door entering the private experimental space, and were informed of and agreed to the nature of the experiment and their rights as a participant, which included a statement that the experiment could be halted at any time.

Participants were then instructed to stand with their toes touching a line on the floor, and were asked to remain there for the duration of the experiment (Figure 5.10). The experimenter then provided instructions about the task to be performed.

Participants were introduced to the robot, and were informed that all of its actions were completely autonomous. Participants were told that the robot would be moving along a straight line throughout the duration of the experiment; a brief demonstration of robot motion was provided, in which the robot autonomously moved back and forth between distances of 3.0 meters and 4.5 meters from the participant, allowing them to familiarize themselves with the robot motion. Participants were told that they would be asked about their preferences regarding the robot's location throughout the experiment.

Participants were then informed that they would be instructing the robot to look at any one of four boxes (of their choosing) located in the room (Figure 5.10), and that they could use speech (in English) and pointing gestures. A vocabulary for robot instructions was provided: for speech, participants were told they could say the words "look at" followed by the name of the shape or color of each box (e.g., "triangle", "circle", "blue", "purple", etc.); for pointing gestures, participants were asked to use their left arm to point to boxes located on their left, and their right arm to point to boxes on their right. This vocabulary was provided to minimize any perceptions the person might have that the robot simply did not understand the words or gestures that they used; thus, the use of the vocabulary attempted to maximize the perception that any failures of the robot were due to other factors.

Participants were told that they would repeat this instruction procedure to the robot many times, and that the robot would indicate whether or not it understood their instructions each time using the head movements and audio cues described in Section 5.4.1.2.

Participants had an opportunity to ask the experimenter any clarifying questions. The experiment proceeded once participant understanding was verified.

5.4.2.2 Phase 2: Pre-interaction Proxemic Preference (*pre*)

The robot autonomously moved to the “far home” location (Figure 5.10). Participants were told that the robot would be approaching them, and to say out loud the word “stop” when the robot reached the ideal location at which the participant would have a *face-to-face conversation*⁹ with the robot. This pre-interaction proxemic preference from the “far home” location is denoted as pre_{far} .

When the participant was ready, the experimenter pressed a PS3 button to start the robot moving. When the participant said “stop”, the experimenter pressed another button to halt robot movement. The experimenter pressed another button to record the distance between the robot and the participant, as measured by the on-board laser.

Once the pre_{far} distance was recorded, the experimenter pressed another button, and the robot autonomously moved to the “near home” location (Figure 5.10); the participant was informed that the robot would be approaching to this location and would stop on its own. The process was repeated with the robot backing away from the participant, and the participant saying “stop” when it reached the ideal location for conversation. This pre-interaction proxemic preference from the “near home” location is denoted as pre_{near} .

Values for pre_{far} and pre_{near} were used to calculate and record the average pre-interaction proxemic preference, denoted as pre ¹⁰.

⁹Related work in human-robot proxemics asks the participant about locations at which they feel *comfortable* (Takayama and Pantofaru, 2009), yielding proxemic preferences very near to the participant. This dissertation focuses on face-to-face human-robot conversational interaction, with proxemic preference farther from the participant (Mead and Matarić, 2014, Torta et al., 2011, 2013), hence the choice of wording.

¹⁰Post hoc analysis revealed no statistically significant difference between pre_{far} and pre_{near} measurements, hence the use of pre .

5.4.2.3 Phase 3: Interaction Scenario

After determining pre-interaction proxemic preferences, the robot returned to the “far home” location. The experimenter then repeated to participants the instructions about the task they were to perform with the robot. When participants verified that they understood the task and indicated that they were ready, the experimenter pressed a button to proceed with the task.

The robot autonomously visited ten pre-determined locations (Figure 5.10). At each location, the robot responded to instructions from the participant to look at one of four boxes located in the room (Figure 5.10). Five instruction-response interactions were performed at each location, after which the robot moved to the next location along its path; thus, each participant experienced a total of 50 instruction-responses interactions. Robot goal locations were in 0.5-meter intervals inclusively between the “near home” location (0.25 meters) and “far home” location (4.75 meters) along a straight-line path in front of the participant (Figure 5.10). Locations were visited in a sequential order; for half of the participants, the robot approached from the “far home” location (i.e., farthest-to-nearest order), and, for the other half of participants, the robot backed away from “near home” location (i.e., nearest-to-farthest order); this was done to reduce any ordering effects (Murata, 1999).

To controllably simulate social signal attenuation at each location, robot performance was artificially manipulated as a function of the distance to the participant—this process is described in detail in Section 5.4.3. After each instruction provided by the participant, the experimenter provided to the robot (via a remote control interface) the ground-truth of the instruction; the robot then determined whether or not it would have understood the instruction based on a prediction from a performance vs. distance curve

(specified by the assigned experimental condition described in Section 5.4.3), and provided either an *affirmative response* or a *negative response* to the participant indicating its successful or failed understanding of the instruction, respectively.

The entire interaction scenario lasted 10-15 minutes, depending on participant speed.

5.4.2.4 Phase 4: Post-interaction Proxemic Preference (*post*)

After the robot visited each of the ten locations, it autonomously returned to the “far home” location. The experimenter then repeated the procedure for determining proxemic preferences described in Section 5.4.2.2. This process generated post-performance proxemic preferences from the “far home” and “near home” locations, as well as their average, denoted $post_{far}$, $post_{near}$, and $post$ ¹¹, respectively.

5.4.2.5 Phase 5: Perceived Peak Location (*perc*)

Finally, after collecting post-interaction proxemic preferences, the experimenter repeated the procedure described in Section 5.4.2.2 to determine participant perceptions of the location of peak performance. This process generated perceived peak performance locations from the “far home” and “near home” locations, as well as their average, denoted $perc_{far}$, $perc_{near}$, and $perc$ ¹², respectively.

This final behavioral measured marked the conclusion of the experiment.

5.4.2.6 Phase 6: Questionnaire

After the experiment concluded, participants were asked to complete a short questionnaire about their experience with the robot. The multi-item questionnaire measured subjective participant perceptions of the robot in terms of the Godspeed questionnaire

¹¹Post hoc analysis revealed no statistically significant difference between $post_{far}$ and $post_{near}$ measurements, hence the use of *post*.

¹²Post hoc analysis revealed no statistically significant difference between $perc_{far}$ and $perc_{near}$ measurements, hence the use of *perc*.

metrics (Bartneck et al., 2009)—*animacy* (9 items; Cronbach’s $\alpha = 0.88$), *anthropomorphism* (6 items; Cronbach’s $\alpha = 0.82$), *competence* (6 items; Cronbach’s $\alpha = 0.84$), *likability* (6 items; Cronbach’s $\alpha = 0.85$), and *safety* (4 items; Cronbach’s $\alpha = 0.55$)—as well as custom metrics of *engagement* (6 items; Cronbach’s $\alpha = 0.76$), *comfort* (6 items; Cronbach’s $\alpha = 0.72$), and *technology adoption* (6 items; Cronbach’s $\alpha = 0.65$). Participants responded to items for each measure on a seven-point Likert scale; items were randomly ordered prior to the study to reduce effects between them, and were consistent across participants.

5.4.3 Experimental Conditions

During each interaction scenario (described in Section 5.4.2.3), robot performance $p(x)$ was artificially varied based on distance x in a between-participants design. Performance values for $p(x)$ were calculated *a priori* for each interaction, and were roughly proportional to a scaled Gaussian distribution parameterized by the following values:

1. μ_x , the “**maximum performance distance**”—the *mean* (center) of the distribution, treated as the *actual location of robot peak performance*;
2. p_{max} , the “**maximum performance value**”—the *maximum* likelihood of recognizing a user instruction;
3. p_{min} , the “**minimum performance value**”—the *minimum* likelihood of recognizing a user instruction; and
4. p_{avg} , the “**average performance value**”—the average likelihood of recognizing a user instruction, indicating how many affirmative and negative responses the robot produced during the interaction (e.g., if $p_{avg} = 0.40$, then the robot provided 20 affirmative responses and 30 negative responses distributed across 50 instructions).

The first three of these parameters— μ_x , p_{max} , and p_{min} —were systematically varied in between-participants conditions, each discussed in Sections 5.4.3.1–5.4.3.3; these manipulations influenced the fourth parameter p_{avg} , which, in turn, impacted the way in which robot responses were distributed across locations between conditions (Table 5.4). In addition to the three **“parameter-varying” conditions**, a uniform distribution of $p(x)$ was also evaluated to serve as a **“baseline” between-participants condition** (discussed in Section 5.4.3.4).

For all conditions, the performance value $p(\mu_x)$ at μ_x was always p_{max} (i.e., the number of affirmative responses was 5 for $p_{max} = 1.00$), and the number of affirmative responses at other distances was always less than that at μ_x to ensure that participants were exposed to an actual singular maximum performance value.

The maximum performance distance was held constant in conditions that varied the maximum and minimum performance values (p_{max} and p_{min} , respectively). A constant distance of $\mu_x = 2.25$ was chosen for three reasons:

1. It was the near-center location between all locations explored, allowing for a balance in the number of affirmative robot responses on each side of the location of maximum performance.
2. It was at a location that is *not* where people initially preferred the robot to be, as Mead and Matarić (2015) reported user proxemic (distance) preferences of ($M = 1.14, SD = 0.49$) meters to the robot used in this study.
3. It was at a location at which the combined performance of actual automated speech and gesture recognition are acceptable and functional (reported in Section 5.3).

An overview of the values for the **baseline** and **parameter-varying** conditions—varying μ_x , p_{max} , and p_{min} (and subsequently p_{avg})—as well as the corresponding distribution of affirmative responses in each condition is presented in Table 5.4. In each condition, the number of affirmative responses was normalized to a desired average robot performance value p_{avg} . The order of the five possible robot responses was randomized at each location.

Each of the four experimental conditions is detailed in corresponding sections below.

		Performance Parameters				Distance, x (meters)									
		μ_x	p_{max}	p_{min}	p_{avg}	0.25	0.75	1.25	1.75	2.25	2.75	3.25	3.75	4.25	4.75
Experimental Condition	BL	N/A	0.40	0.40	0.40	2	2	2	2	2	2	2	2	2	2
	Vary μ_x	0.75	1.00	0.00	0.40	4	5	4	3	1	1	1	1	0	0
		1.25	1.00	0.00	0.40	3	4	5	4	3	1	0	0	0	0
		1.75	1.00	0.00	0.40	1	3	4	5	4	3	0	0	0	0
		2.25	1.00	0.00	0.40	0	1	3	4	5	4	3	0	0	0
		2.75	1.00	0.00	0.40	0	0	0	3	4	5	4	3	1	0
		3.25	1.00	0.00	0.40	0	0	0	0	3	4	5	4	3	1
		3.75	1.00	0.00	0.40	0	0	0	0	1	3	4	5	4	3
		4.25	1.00	0.00	0.40	0	0	1	1	1	1	3	4	5	4
	Vary p_{max}	2.25	0.20	0.00	0.02	0	0	0	0	1	0	0	0	0	0
		2.25	0.40	0.00	0.20	1	1	1	1	2	1	1	1	1	0
		2.25	0.60	0.00	0.30	1	1	2	2	3	2	2	1	1	0
		2.25	0.80	0.00	0.40	1	1	3	3	4	3	3	1	1	0
	Vary p_{min}	2.25	1.00	0.20	0.50	1	1	3	4	5	4	3	2	1	1
		2.25	1.00	0.40	0.60	2	2	3	4	5	4	3	3	2	2
		2.25	1.00	0.60	0.70	3	3	3	4	5	4	4	3	3	3
2.25		1.00	0.80	0.82	4	4	4	4	5	4	4	4	4	4	
# Affirmative Responses vs. Distance															

Table 5.4: The distribution of affirmative responses provided by the robot across **baseline (BL)** and **parameter-varying conditions**; note that average performance values p_{avg} vary as well. Manipulated values in each condition are highlighted in **bold italics**.

5.4.3.1 Condition: Maximum Performance Distance (μ_x)

To explore the space of human responses to robot performance differences at a variety of distances, the maximum performance distance μ_x was varied by selecting the eight locations non-inclusively between the “near home” and “far home” locations (Figures 5.10 and 5.12]; the “near home” and “far home” locations were not included in the set of μ_x to ensure that participants were always exposed to an actual *peak* in performance, rather than just a *trend*. The value of μ_x was varied between participants. The number of affirmative responses was normalized to 20 (40%) to ensure a consistent user experience of average robot performance $p_{avg} = 0.40$ for different values of μ_x .

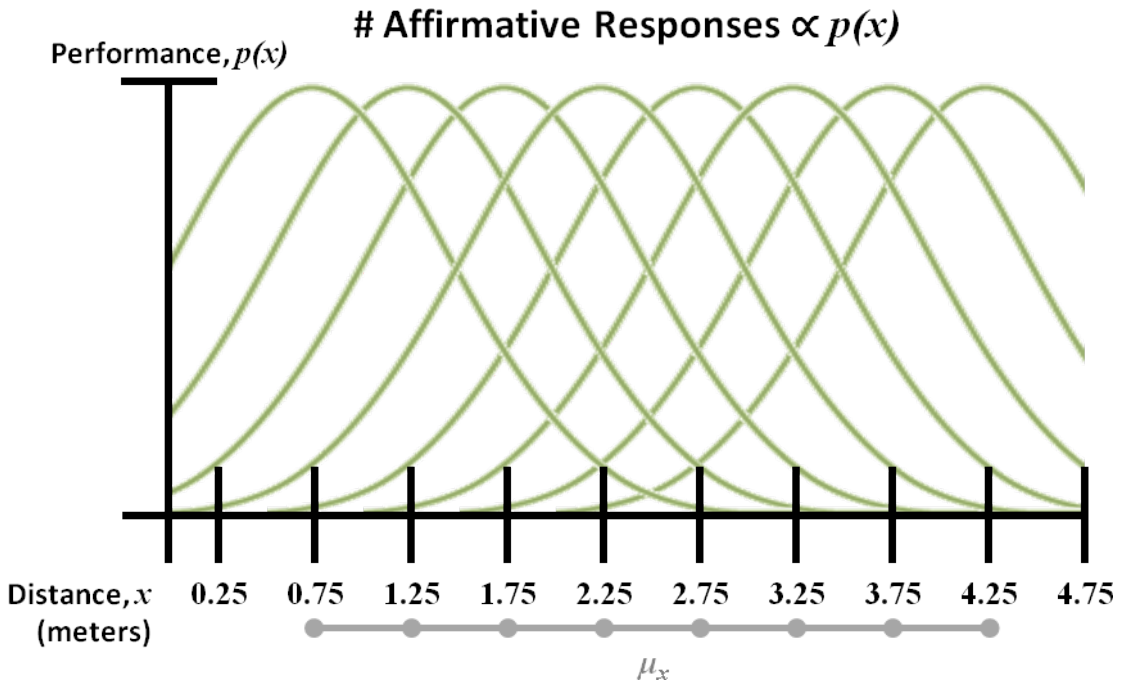


Figure 5.12: Manipulation condition varying the **maximum performance distance**, $\mu_x = \{0.75, 1.25, 1.75, 2.25, 2.75, 3.25, 3.75, 4.25\}$ meters.

5.4.3.2 Condition: Maximum Performance Value (p_{max})

To investigate how people responded to differences in robot maximum performance, p_{max} was varied at a constant location of maximum performance, $\mu_x = 2.25$. Values of p_{max} were 0.20, 0.40, 0.60, and 0.80 (Figure 5.13); for analysis, $p_{max} = 1.00$ at $u_x = 2.25$ ($p_{avg} = 0.40$) collected in the condition varying μ_x (Section 5.4.3.1) was also considered as part of this group. Average performance values naturally varied in this condition, with $p_{avg} = \{0.02, 0.20, 0.30, 0.40\}$, respectively.

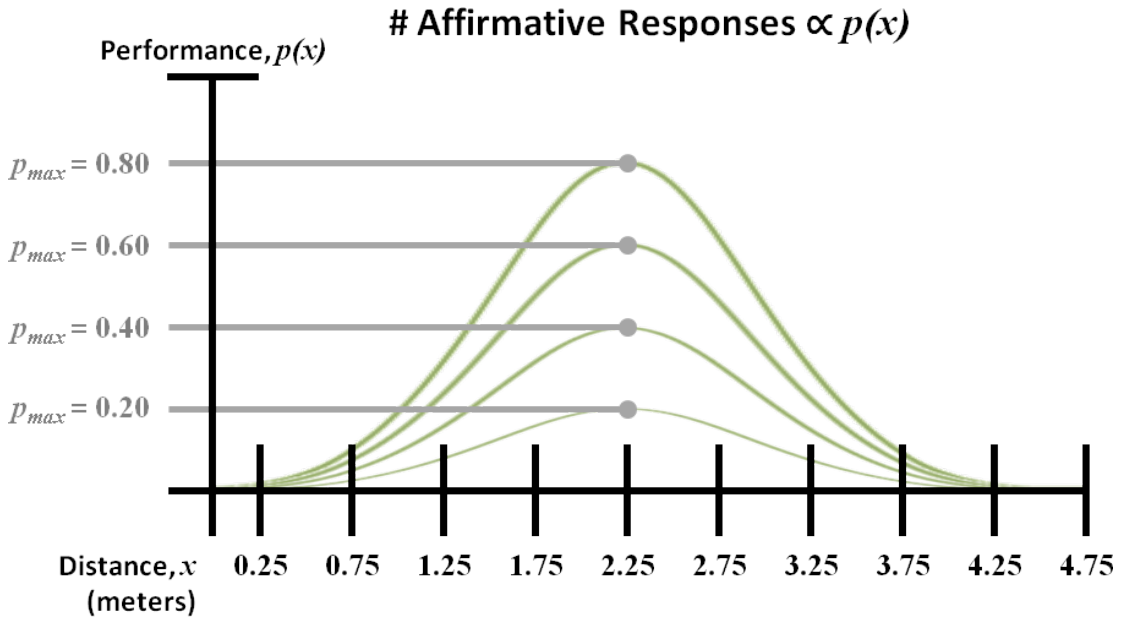


Figure 5.13: Manipulation condition varying the **maximum performance value**, $p_{max} = \{0.20, 0.40, 0.60, 0.80\}$.

5.4.3.3 Condition: Minimum Performance Value (p_{min})

The experiment also investigated how people responded to differences in robot minimum performance. In this condition, p_{min} was varied in a way similar to that of p_{max} (Section 5.4.3.2), with $\mu_x = 2.25$. Values of p_{min} were 0.20, 0.40, 0.60, and 0.80; for analysis, $p_{min} = 0.00$ at $u_x = 2.25$ ($p_{avg} = 0.40$) collected in the condition varying μ_x (Section 5.4.3.1) was also considered as part of this group. Average performance values naturally varied in this condition, with $p_{avg} = \{0.50, 0.60, 0.70, 0.82\}$, respectively.

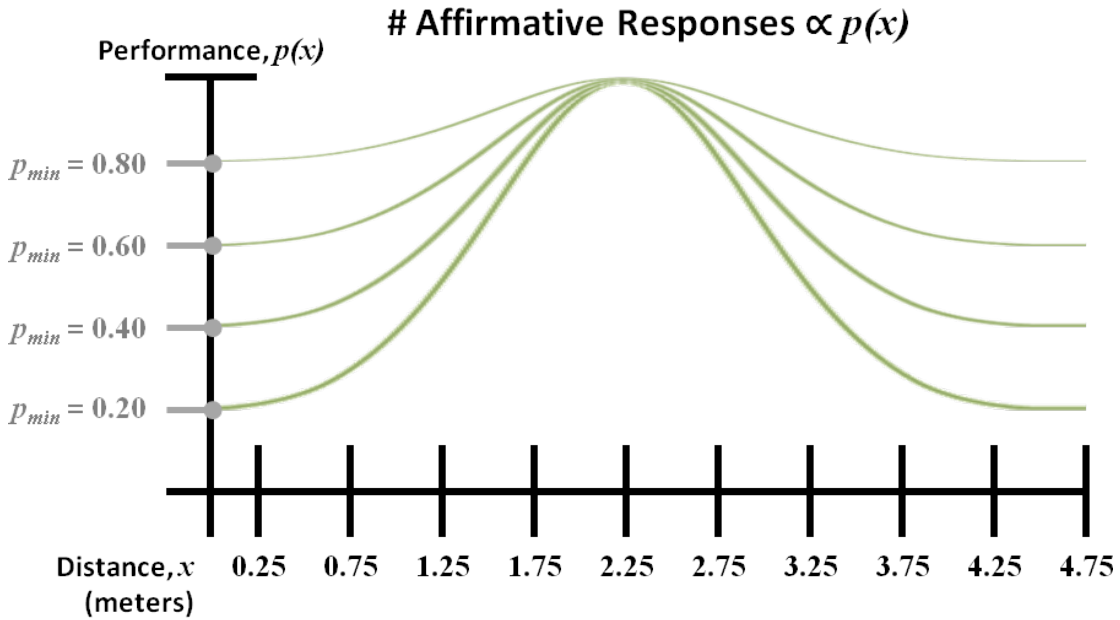


Figure 5.14: Manipulation condition varying the **minimum performance value**, $p_{min} = \{0.20, 0.40, 0.60, 0.80\}$.

5.4.3.4 Condition: Baseline ($p(x) = p_{max} = p_{min}$)

In this condition, robot performance was the same ($p_{max} = p_{min} = 40\%$) at all locations (Figures 5.15 and Table 5.4). Thus, at each of the ten locations visited, the robot provided two affirmative and three negative responses, respectively. This condition served as a baseline of participant proxemic preferences within the task.

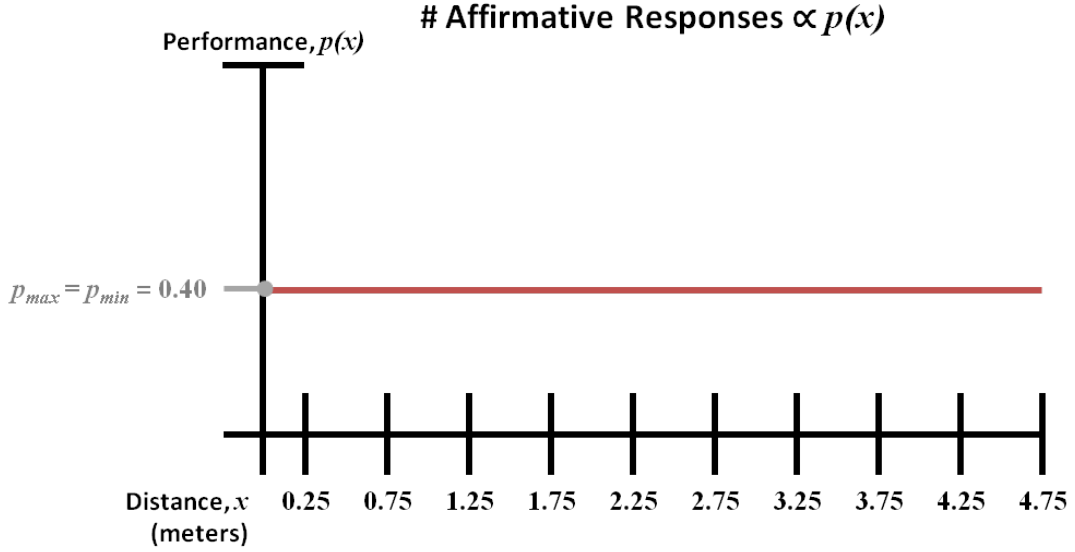


Figure 5.15: Baseline condition in which performance vs. distance is represented by a uniform distribution, $p(x) = p_{max} = p_{min} = 0.40$.

5.4.4 Experimental Hypotheses

Within these conditions, the study explored three central hypotheses:

H1: In the **baseline condition**, there will be no significant change in participant proxemic preferences.

H2: In the μ_x -**varying condition**, participants will be able to identify a relationship between robot performance and human-robot proxemics.

H3: In the μ_x -**varying condition**, participants will adapt their proxemic preferences to improve robot performance.

5.4.5 Participants

A total of 180 participants (90 male, 90 female) were recruited from campus community at the University of Southern California. Participant race was diverse (105 white/Caucasian, 57 Asian, 8 Latino/Latina, 7 black/African-American, and 3 mixed-race). All participants reported proficiency in English and had lived in the United States for at least two years (i.e., had acclimated to U.S. culture). The average age (in years) of participants was 21.34 ($SD = 4.26$), ranging from 18 to 39. Based on a seven-point scale, participants reported moderate familiarity with technology ($M = 4.04$, $SD = 0.93$). The average participant height (in meters) was 1.72 ($SD = 0.14$), ranging from 1.49 to 1.96. Related work reports how human-robot proxemics is influenced by participant gender and technology familiarity (Takayama and Pantofaru, 2009), culture (Eresha et al., 2013), and height (Hiroi and Ito, 2011, Rae et al., 2005).

The 180 participants were randomly assigned to a condition, with $N = 20$ in the **baseline condition** $N_{\mu_x} = 80$, $N_{p_{max}} = 40$, $N_{p_{min}} = 40$ in the conditions varying maximum performance distance (μ_x), maximum performance value (p_{max}), and minimum performance value (p_{min}), respectively. Within each **parameter-varying condition**, the participants were randomly assigned to one of the subconditions (e.g., one of the eight maximum performance distances; described in Section 5.4.3) with $N = 10$ for each subcondition. Neither the participant nor the experimenter was aware of the assigned condition.

5.4.6 Data Analysis and Results: Behavioral Measures

The data collected were analyzed to test the three hypotheses (described in Section 5.4.4) with respect to pre-interaction (*pre*), post-interaction (*post*), and perceived peak location (*perc*) *behavioral* measures. Each hypothesis is evaluated in the corresponding sections below, followed by a discussion of the implications for autonomous sociable robots and human-robot proxemics.

To provide a baseline of the robot used for comparison in general human-robot proxemics, pre-interaction proxemic preferences (*pre*) across all conditions ($N = 180$) were consolidated and analyzed, as the data had not yet been influenced by the manipulation in the assigned experimental condition. The participant pre-interaction proxemic preference (in meters) was determined to be 1.14 ($SD = 0.49$) for the Bandit robot system, which is consistent with (Mumm and Mutlu, 2011) and (Mead and Matarić, 2014), but twice as far away as related work has reported for robots of a similar form factor (Takayama and Pantofaru, 2009, Walters et al., 2009).

5.4.6.1 H1: Pre- vs. Post-interaction Locations

To test **H1**, average pre-interaction proxemic preferences (*pre*) were compared to average post-interaction proxemic preferences (*post*) of participants in the **baseline condition**.

A paired *t*-test revealed a statistically significant change in participant proxemic preferences between *pre* ($M = 1.12, SD = 0.51$) and *post* ($M = 1.39, SD = 0.63$); $t(38) = 1.49, p = 0.02$. Thus, hypothesis **H1** is rejected.

This result suggests that there might be something about the context of the interaction scenario itself that influenced participant proxemic preferences. To address any influence the interaction scenario might have on subsequent analyses, a *contextual offset* θ was defined as the average difference in participant post-interaction and pre-interaction proxemic preferences ($M = 0.27, SD = 0.48$); this θ value will be subtracted from (*post* – *pre*) values in Section 5.4.6.3 to normalize for the interaction context.

5.4.6.2 H2: Perceived vs. Actual Peak Locations

To test **H2**, participant perceived locations of peak performance (*perc*) were compared to actual locations of peak performance (*peak* = μ_x) in the μ_x -**varying conditions** (Figure 5.16).

Steven’s Power Law, ax^b , has previously been used to model human distance estimation as a function of actual distance (Murata, 1999), and is generally well representative of human-perceived vs. actual stimuli (Stevens, 2007). However, existing Power Laws relevant to the dissertation only pertain to distances of 3–23 meters, which are beyond the range of the natural face-to-face communication with which this work is concerned. Thus, the goal here is to model the collected data to establish a Power Law for *perc* vs. *peak* at locations more relevant to HRI (0.75–4.25 meters), which can then be evaluated to test **H2**.

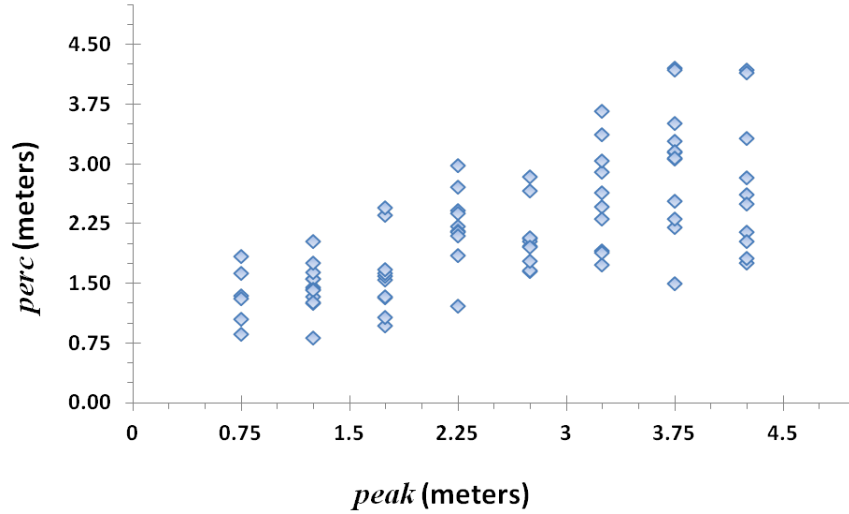


Figure 5.16: Participant perceived location of robot peak performance (*perc*) vs. actual location of robot peak performance (*peak*). Note the heteroscedasticity of the data, which prevents us from performing traditional statistical analyses without first transforming the data (shown in Figure 5.17).

Immediate observations of the dataset suggested that the data appear to be heteroscedastic (Figure 5.16)—in this case, the variance seems to increase with distance from the participant, which means traditional statistical tests should not be used. The Breusch-Pagan test for non-constant variance (NCV) confirmed this intuition; $\chi^2(1, N = 100) = 15.79, p < 0.001$. A commonly used and accepted approach to alleviate heteroscedasticity is to transform the *perc* and *peak* data to a log-log scale. While not applicable to all datasets, this approach served as an adequate approximation for the collected data (Figure 5.17); it also supports the application of a regression analysis to determine parameter values for the Power Law coefficient and exponent, $a = 1.3224$ and $b = 0.5132$, respectively. With these parameters, it was revealed that *peak* was a strongly correlated and very significant predictor of *perc*; $R^2 = 0.4951, F(1, 78) = 76.48, p < 0.001$. Thus, hypothesis **H2** is supported.

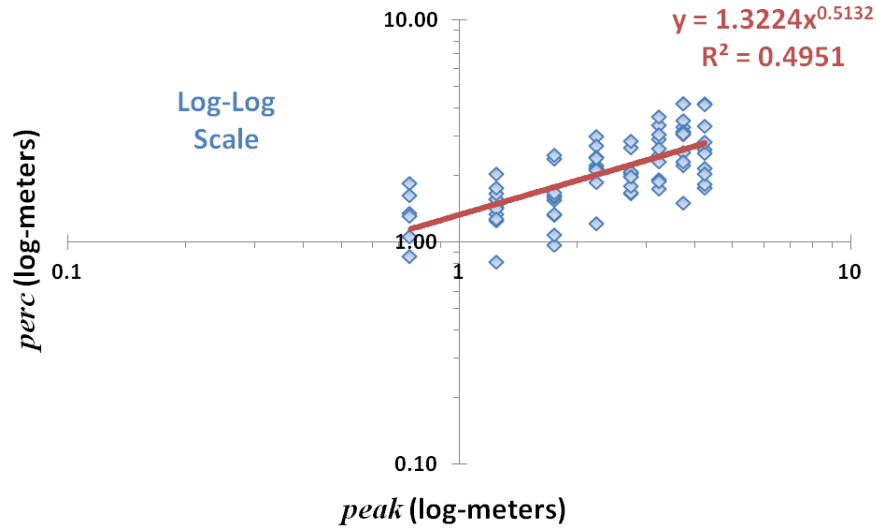


Figure 5.17: Participant perceived location of robot peak performance (*perc*) vs. actual location of robot peak performance (*peak*) on a log-log scale, reducing the effects of heteroscedasticity and allowing us to perform regression to determine parameters of the Power Law, ax^b .

This result suggests that people are able to identify a relationship between robot performance and human-robot proxemics, but they will predictably underestimate the distance, x , to the location of peak performance based on the Power Law equation $1.3224x^{0.5132}$. While human estimation of the location of peak performance is suboptimal, it is possible that repeated exposure to the robot over multiple sessions might yield more accurate results.

5.4.6.3 H3: Preferences vs. Peak Locations

To test **H3**, changes in participant pre-/post-interaction proxemic preferences ($post - pre - \theta$) were compared to the distance from the participant pre-interaction proxemic preference to either a) the actual location of robot peak performance ($peak - pre$) (Figure 5.18), or b) the perceived location of robot peak performance ($perc - pre$) (Figure 5.19), both in the μ_x -**varying conditions**.

Data for ($post - pre - \theta$) vs. both ($peak - pre$) and ($perc - pre$) were heteroscedastic, as indicated by Breusch-Pagan NCV tests: $\chi^2(1, N = 80) = 18.81, p < 0.001$; and $\chi^2(1, N = 80) = 13.55, p < 0.001$; respectively. This is intuitive, as the data for perceived ($perc$) vs. actual ($peak$) locations of peak performance were also heteroscedastic (Figure 5.16). The log-transformation approach used in Section 5.4.6.2 did not perform well in modeling these data; thus, an alternative approach was needed. A Generalized Linear Model (Nelder and Wedderburn, 1972) was utilized because it models the variance of each measurement separately as a function of predicted values and allows appropriate statistical tests for significance to be applied.

Changes were modeled in participant proxemic preferences ($post - pre - \theta$) vs. distance from pre-interaction proxemic preference to the actual location of peak performance ($peak - pre$). In the ideal situation (for the robot), these match one-to-one—in other words, the participant meets the needs of the robot entirely by changing proxemic preferences to be centered at the peak of robot performance. Unfortunately for the robot, this was not the case. A strongly correlated and statistically significant relationship was detected between participant proxemic preference change and distance from pre-interaction preference to the peak location ($R^2 = 0.5474, \beta = 0.5361, t(78) = 9.71, p < 0.001$), but participant preference change was predictably poor ($\beta = 0.5361$) with respect to robot location of peak performance (Figure 5.18).

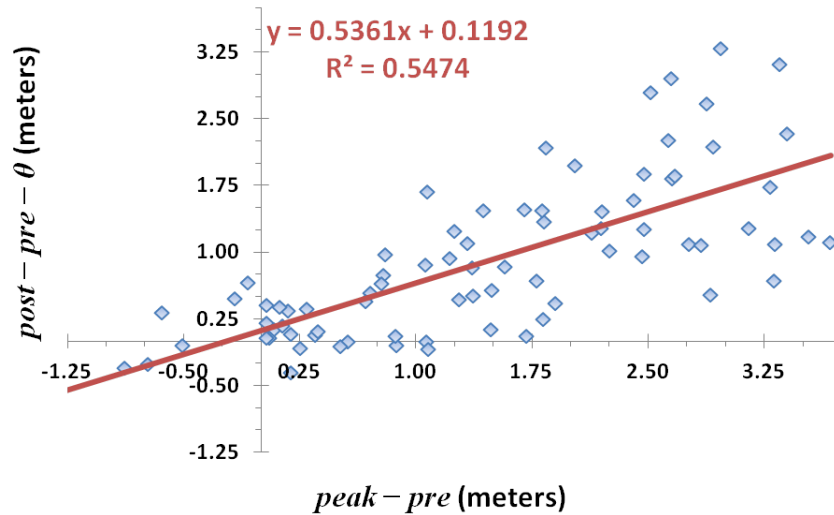


Figure 5.18: Changes in participant pre-/post-interaction proxemic preferences (pre and $post$, respectively; θ is the *contextual offset* defined in Section 5.4.6.1) vs. distance from participant pre-interaction proxemic preference (pre) to the actual location of robot peak performance ($peak$).

Recall that results reported for **H2** (Section 5.4.6.2) suggested that, while people do perceive a relationship between robot performance and distance, their ability to accurately identify the location of robot peak performance diminishes based on the distance to it as governed by a Power Law. Were participants *trying* to maximize robot performance, but simply adapting their preferences to a suboptimal location?

This question was investigated by considering changes in participant proxemic preferences ($post - pre - \theta$) vs. distance from pre-interaction proxemic preference to the perceived location of peak performance ($perc-pre$). If the participant was adapting their proxemic preferences to accommodate the needs of the robot, then these should match one-to-one. A Generalized Linear Model was fit to these data, and yielded a strongly correlated and statistically significant relationship between changes in proxemic preferences and perceptions of robot performance ($R^2 = 0.5421, \beta = 0.9275, t(78) = 9.61, p < 0.001$) (Figure 5.19). Thus, hypothesis **H3** is supported.

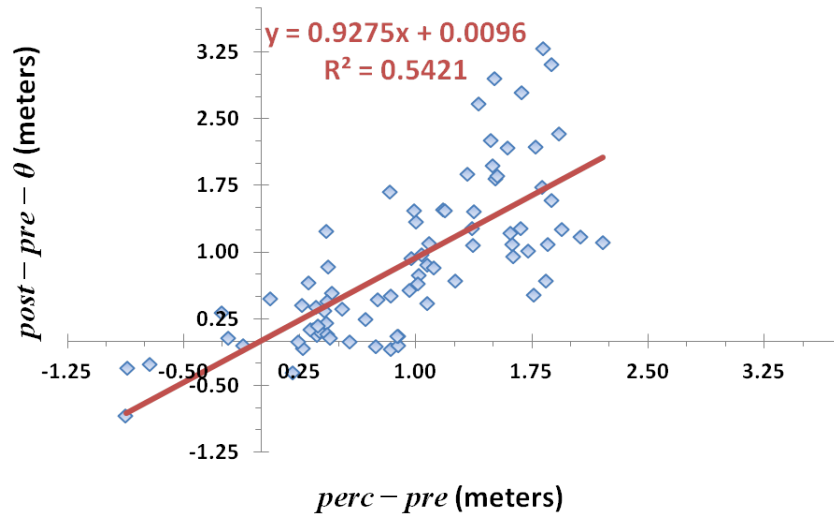


Figure 5.19: Changes in participant pre-/post-interaction proxemic preferences (pre and $post$, respectively; θ is the *contextual offset* defined in Section 5.4.6.1) vs. distance from participant pre-interaction proxemic preference (pre) to the perceived location of robot peak performance ($perc$).

The near one-to-one relationship ($\beta = 0.9275$) between post-interaction proxemic preferences and participant perceptions of robot peak performance is compelling, suggesting that participants adapted their proxemic preferences almost entirely to improve robot performance (an *objective* measure) in the interaction.

5.4.7 Data Analysis and Results: Subjective Measures

Bidirectional stepwise multivariate linear regression was used to model and analyze the relationships between participant subjective experiences of the robot (described in Section 5.4.2.6) and the four parameters of the distribution of robot responses (i.e., μ_x , p_{max} , p_{min} , and p_{avg} ; described in Section 5.4.3) across all conditions. Of the eight subjective factors considered, significant effects were found for five of them—*competence*, *anthropomorphism*, *engagement*, *likability*, and *technology adoption*. Figure 5.20 illustrates the correlation coefficient (β) and statistical significance (p) of each predictor-measure pair represented in the resulting regression models. For these analyses, t -tests determined statistical significance (p) and adjusted R^2 values (R_A^2) indicate the explanatory power of the model based on the number of predictors. A detailed summary of the results is provided below and further discussed in Section 5.4.8.2.

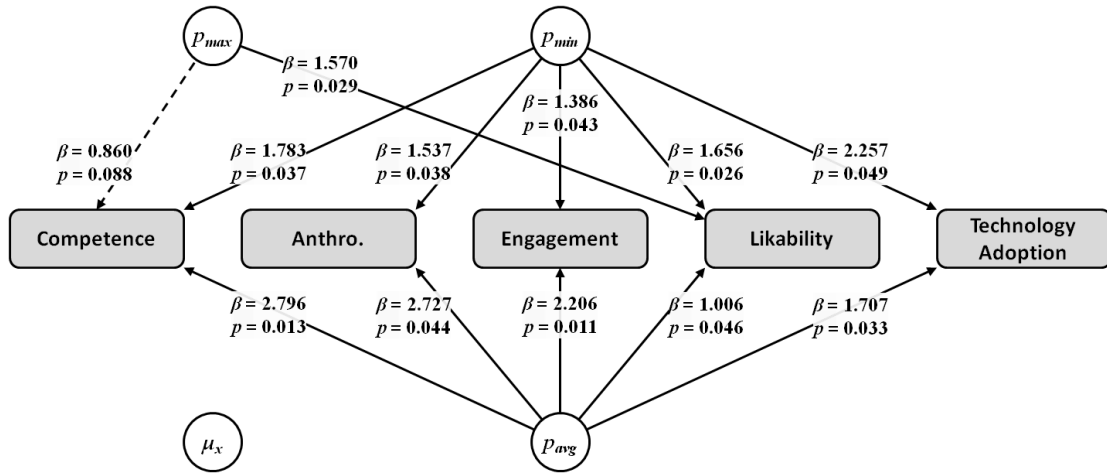


Figure 5.20: The significant relationships modeled between the manipulated predictor variables and subjective measures. The correlation coefficient (β) and statistical significance (p) for each predictor-measure pair is presented along the connecting line; a dotted line indicates marginal significance.

5.4.7.1 Subjective Measure: Competence

Values of average performance (p_{avg}) and minimum performance (p_{min}) were found to be significant predictors of perceived competence ($R_A^2 = 0.425, \beta_0 = 1.352$), with ($\beta = 2.796, p = 0.013$) and ($\beta = 1.783, p = 0.037$), respectively; furthermore, maximum performance (p_{max}) was a marginal predictor of perceived competence ($\beta = 0.860, p = 0.088$). This indicates that robot performance factors should be strongly considered when deploying a sociable robot, as the perceptions of intelligence attributed by a co-present human user will dramatically vary based on how well the robot functions. Average performance should be as high as possible, as it is the deciding factor in perceptions of intelligence. Situations in which the robot fails should be avoided, as they hinder perceived robot competence. High robot maximum performance should be a tertiary objective, as it is only marginally predictive of perceived intelligence.

5.4.7.2 Subjective Measure: Anthropomorphism

Values of average performance (p_{avg}) and minimum performance (p_{min}) were found to be significant predictors of perceived anthropomorphism ($R_A^2 = 0.412, \beta_0 = 1.679$), with ($\beta = 2.727, p = 0.044$) and ($\beta = 1.537, p = 0.038$), respectively. This suggests that a robot that is consistently performing well at many locations will receive the greatest attributions of agency and human-likeness; furthermore, reducing the number of times in which the robot fails will further increase these attributions.

5.4.7.3 Subjective Measure: Engagement

Values of average performance (p_{avg}) and minimum performance (p_{min}) were found to be significant predictors of participant engagement ($R_A^2 = 0.378, \beta_0 = 2.878$), with ($\beta = 2.206, p = 0.011$) and ($\beta = 1.386, p = 0.043$), respectively. Thus, for increased user engagement, the robot should perform well on average and should avoid situations in which it would likely fail. In situations in which the robot was consistently failing (i.e., when both average and minimum performance were low), it was clear from experimenter observations that participants were significantly less engaged than when the robot was performing better. Interestingly, experimenter observations of participants also suggested that participants appeared bored (i.e., less engaged) when minimum performance was at the highest level explored ($p_{min} = 0.80$), as the robot performed very well at every location, though participant questionnaire responses did not support this observation. Related work (Short et al., 2010) has shown that people are more engaged by a robot that produces imperfect or unexpected behaviors (in Short et al. (2010), the robot cheated during a game of rock-paper-scissors), which might offer insights into this observation.

5.4.7.4 Subjective Measure: Likability

Values of minimum performance (p_{min}), maximum performance (p_{max}), and average performance (p_{avg}) were found to be significant predictors of perceived likability ($R_A^2 = 0.317, \beta_0 = 1.634$), with ($\beta = 1.656, p = 0.026$), ($\beta = 1.570, p = 0.029$), and ($\beta = 1.066, p = 0.046$), respectively. People do not like it when the robot consistently fails—this is intuitive, but has not been well quantified in the field. Increased average performance provided some support, and could be a deciding factor in a user’s liking of the robot.

5.4.7.5 Subjective Measure: Technology Adoption

Values of minimum performance (p_{min}) and average performance (p_{avg}) were found to be significant predictors of participant technology adoption ($R_A^2 = 0.322, \beta_0 = 2.397$), with ($\beta = 2.257, p = 0.049$) and ($\beta = 1.707, p = 0.033$), respectively. Repetitive failures can dramatically reduce the possibility of a user adopting the technology, even if the robot performs well on average. Maximum performance is therefore not the most important factor in design; rather, avoidance of failure is more important.

5.4.8 Discussion

The results obtained from this study have significant implications for the design of social robots and autonomous robot proxemic control systems. Sections 5.4.8.1 and 5.4.8.2 discuss these implications with respect to *behavioral* and *subjective* results, respectively.

5.4.8.1 Behavioral Results

The primary implication of the behavioral results reported in Section 5.4.6 is that the proxemic preferences of human users will change over time as the user interacts with and comes to understand the needs of the robot, and these changes will improve robot performance (an *objective* measure of automated speech and gesture recognition rates).

As illustrated in the data-driven approach of proxemics and multimodal communication (Section 4.3.2), the locations of on-board sensors for social signal recognition (e.g., microphones and cameras), as well as the automated speech and gesture recognition software used, can have significant impacts on the performance of the robot in autonomous face-to-face social interactions. The behavioral results reported for this study suggest that people will adapt their behavior in an effort to improve robot performance, so it is anticipated that human-robot proxemics will vary between robot platforms with different hardware and software configurations based on factors that are:

1. not specific to the user (unlike culture (Eresha et al., 2013), or gender, personality, or familiarity with technology (Takayama and Pantofaru, 2009));
2. not observable by the user (unlike height (Hiroi and Ito, 2011, Rae et al., 2005), amount of eye contact (Mumm and Mutlu, 2011, Takayama and Pantofaru, 2009), or vocal parameters (Walters et al., 2008)); and
3. not observable by the robot developer.

User understanding of the relationship between robot performance and human-robot proxemics is a latent factor that only develops through repeated interactions with the robot (perhaps expedited by the robot communicating its predicted error); fortunately, the results presented here indicate that user understanding will develop in a predictable way. Thus, it is recommended that developers of sociable robots consider and perhaps model robot performance as a function of conditions that might occur in dynamic proxemic interactions with human users to better predict and accommodate how the people will actually use the technology. This dynamic relationship, in turn, will enable more rich autonomy for sociable robots by improving the performance of their own automated recognition systems.

If developers adopt models of robot performance as a factor contributing to human-robot proxemics, then it follows that proxemic control systems should be designed to expedite the process of autonomously positioning the robot at an optimal distance from the user to maximize robot performance while still accommodating the initial personal space preferences of the user. This is the focus of the data-driven approach in this dissertation—treating proxemics as an optimization problem that considers the production and perception of social signals (speech and gesture) as a function of distance and orientation. An objective of this study was to address questions regarding whether or not users would accept a robot that positions itself in locations that might differ from

their initial proxemic preferences. The behavioral results reported (specifically, in Section 5.4.6.3) support the notion that user proxemic preferences will change through interactions with the robot as its performance is observed, and that the new user proxemic preference will be at the *perceived* location of robot peak performance. An extension of this result is that, through repeated interactions, user proxemic preferences will further adapt and eventually converge to the *actual* location of robot peak performance.

5.4.8.2 Subjective Results

The primary implication of the subjective results reported in Section 5.4.7 is that the relationship between robot performance and participant subjective experiences should be strongly considered by developers if sociable robots are to be successfully deployed.

Overall, the subjective results of this study suggest that *average performance* across multiple interactions, rather than *maximum performance* in a select few interactions, should be a focus of robot developers, especially when considering human perceptions of robot competence, anthropomorphism, and engagement. Furthermore, in consideration of all reported significant effects, a robot developer should have the goal of minimizing the amount of user exposure to robot failure (as indicated by the effect of *minimum performance* values); if the robot can predict conditions in which it will most likely fail at recognizing user input (e.g., at certain physical locations, or based on the detected sound pressure level of speech or the regions of space occupied by gestures (Mead and Matarić, 2014)), then the robot could autonomously determine how to react in the situation. This is addressed by the data-driven approach presented in this dissertation—the robot selects proxemic behaviors that it predicts will maximize its ability to autonomously recognize human multimodal communication (speech and gesture), thus attempting to optimize maximum and average system performance.

While no statistically significant relationship was identified between the location of maximum performance (μ_x) and participant subjective experiences along any measure, this serves as an important result in human-robot proxemics. Previous work in human-robot proxemics predicts changes in factors such as perceived safety (Henkel et al., 2012), comfort/intimacy (Henkel et al., 2012, Mumm and Mutlu, 2011, Takayama and Pantofaru, 2009), and likability (Henkel et al., 2012, Mumm and Mutlu, 2011); similar changes are also predicted in human-human proxemics literature (Argyle and Dean, 1965, Hall, 1966). This could be explained by the factors that govern proxemic behavior and how people perceive them. Hall’s (Hall, 1963) psychophysical representation suggests that human-human proxemic behavior can be represented by the sensory experiences of people at different distances, and that the four psychological proxemic zones—public, social, personal, and intimate—that encode the interpersonal relationship between two people (Hall, 1966) are characterized by these sensory experiences. The implication then is that if certain sensory experiences are different or absent during a face-to-face interaction between two agents (human or robot), then perceptions of factors such as comfort or intimacy could not be predicted by distance alone—in short, sensory experience (of multimodal communication), rather than distance, is the significant variable in predicting these subjective measures (as suggested by the results reported in this study). This adds strong support to the psychophysical representation of proxemic behavior, which is at the core of the computational framework presented in this dissertation (Section 3).

5.5 Summary

This chapter investigated the deployment of the computational framework of situated proxemics and multimodal communication in autonomous systems for sociable robots. These systems were implemented for automated human behavior recognition (using the heuristic approach) and autonomous robot behavior generation (using the data-driven approach), and were evaluated in three studies based on objective measure (system performance), subjective measures (user perceptions), and behavioral measures (changes in user behavior). The implemented systems are publicly available as open-source software; the data collected in the studies performed are available upon request. The results of these studies have significant implications for the design of sociable robots, as well as autonomous proxemic behavior and multimodal communication systems for human-robot interaction.

The next and final chapter provides an overview of the work performed for this dissertation.

Chapter 6

Summary and Conclusions

6.1 Contributions

This dissertation presented a novel representation and computational framework that relates proxemics and multimodal communication in both human-human and human-robot interactions. The framework is grounded in social science literature, and computational models implementing the framework were developed from human-robot pose, speech, and gesture interaction data, which are available upon request. Applications of the approach were discussed within the contexts of automated human behavior recognition and autonomous robot control systems, which are publicly available as open-source software. These systems were evaluated in human-robot interaction scenarios, and have significant implications for the design of robot social behaviors in face-to-face HRI. This work establishes a foundational component of HRI, and contributes to the understanding of the underlying processes that govern proxemic behavior and multimodal communication in both human-human and human-robot interactions.

The following are the primary contributions of this work:

1. *An extensible unifying framework for situated proxemics and multimodal communication for both human-human and human-robot interactions.* The framework considers how both humans and robots experience social signals in face-to-face interactions. Data collections were conducted to inform probabilistic graphical models that predict how speech and gesture are produced (transmitted) and perceived (received) by both humans and robots at different distances and under environmental interference.
2. *Proxemic feature extraction and behavior recognition systems.* The system automatically extracts proxemic features based on three feature representations—physical, psychological, and psychophysical. These features are used to recognize transitions into (initiation) and out of (termination) co-present social interactions. A comparison of representations is provided, in which the psychophysical representation of proxemic behavior outperforms traditional physical and psychological representations, suggesting a more powerful approach to the recognition of spatiotemporal interaction cues in both human-human and human-robot interactions.
3. *Probabilistic models for autonomous robot generation of proxemic behavior and multimodal communication.* Using the computational framework in this dissertation, proxemic and communication behavior are unified in a probabilistic graphical model that represents the production and perception of social signals (e.g., speech and gesture) as a function of interagent pose (distance and orientation). Pose, speech, and gesture parameters are selected to maximize social signal recognition rates for all agents (human and robot) in the interaction.

4. *A method for adapting robot proxemic and multimodal communication parameters in complex environments.* The approach has implications for situated robot behavior generation in complex environments (in which there are loud noises or visual occlusions), and for technology personalization in complex interactions with a focus on socially assistive contexts with people with special needs, such as those with hearing or visual impairments or sensitivities.
5. *Objective, subjective, and behavioral evaluations of proxemic and multimodal communication control systems.* Experiments were conducted that a) demonstrated accurate predications of human social signal production and robot social signal recognition (objective measures), b) related user perceptions of the robot to its ability to recognize social signals (subjective measures), and c) demonstrated that human users naturally adapt their own proxemic preferences to improve robot performance in social contexts (behavioral measures).

This dissertation also provides software and a corpus of public human-human and human-robot interaction data for use by researchers worldwide, serving to inform, validate, and extend longstanding research in the both HRI and the social sciences.

The following are the secondary contributions of this work:

1. *Implemented and validated open-source software systems.* The automated human behavior recognition and autonomous robot behavior control systems are publicly available in the Social Behavior Library (SBL), an open-source software suite that provides generic computational models of social behavior for HRI.
2. *Large corpora of human-human and human-robot interaction data.* This dissertation provides uniquely comprehensive corpora of human-human and human-robot interaction data, which are publicly available upon request to facilitate research in both robotics and social science communities.

6.2 Future Work: Adaptation in Complex Interactions

In Section 4.3, data collections were conducted to inform the parameters of the Bayesian network in the data-driven approach (Mead and Matarić, 2014). In Section 5.3, the resulting models were implemented in an autonomous proxemic behavior and multimodal communication control system for a sociable robot (Mead and Matarić, 2012). These controllers used models that are based on an averaging of results *across* participants; however, the models do not account for variability *between* participants, which might result in inadequate robot behavior (e.g., the robot speaks too softly or loudly). Implications of the across vs. between evaluations, as well as methods of adapting models of human social signal production and perception, are discussed below.

Across participants (i.e., considered as a collective), a strong positive linear relationship between human speech output levels and distance to both human and robot interaction partners was identified, as discussed in Section 4.3.2.2. No statistically significant differences in human gesture output levels at different distances were observed.

Between participants (i.e., considered as individuals), significant variability was identified in human speech and gesture output levels (e.g., some participants spoke louder or occupied more space when gesturing)—and this variability has considerable impact on robot speech and gesture input levels (automated recognition rates). Similarly, human users might vary in their own ability to perceive speech and gestures (e.g., if they have hearing or visual impairments or sensitivities); the models of robot speech and gesture output presented are static and do not accommodate for these differences.

These insights led to the consideration of methods for dynamic adaptation of the data-driven models of human social signal production and perception attenuated by interaction pose (distance and orientation). As speech varied the most with distance and between users, adaptation to human speech output and input levels (SOL_{HR} and SIL_{HR} , respectively) should be investigated to extend the models provided by this dissertation. Specifically, the research would answer the questions:

1. How can the robot dynamically adapt its models of *human speech levels* (SOL_{HR}) to select better interagent poses for understanding user social signals?
2. How can the robot dynamically adapt its models of *human speech input levels* (SIL_{HR}) to better select its own multimodal communication parameters for user understanding of its social signals?

These questions are characterized by distinct implications and challenges—subsequently, each warrants its own unique solution. Approaches to the adaptation these models of human speech output and input levels are proposed below in Sections 6.2.1 and 6.2.2, respectively.

6.2.1 Adaptive Models of Human Speech Output Levels (SOL_{HR})

In the implementation of the data-driven approach (Section 5.3), the robot selects and moves to a pose that will most likely maximize its perceptions of social signals (including SIL_{RH}) based on models of human social signal production (including SOL_{HR}) (Section 5.3.1.1; Figure 5.4); thus, the robot will not perform as well if the models of human speech output are inaccurate.

As shown in Figure 4.13, human speech output levels (SOL_{HR}) can be estimated by a linear function of interagent pose (POS) (Mead and Matarić, 2014):

$$SOL_{HR}(POS) = b + m \times POS + \epsilon \quad (6.1)$$

where b is the intercept, m is the slope, and ϵ is the residual standard error (RSE) of the line, each of which the robot must adapt.

Given a speech input level (SIL_{RH}) measured by the robot's microphones (in dB SPL) at an interagent pose (POS), the objective is to estimate the actual human speech output level (SOL_{HR}), defined by an inverse attenuation function:

$$A^{-1}(SIL_{RH}) = A_{\beta}^{-1}(A_d^{-1}(A_{\alpha}^{-1}(SIL_{RH}))) \quad (6.2)$$

where A_{β}^{-1} , A_d^{-1} , and A_{α}^{-1} are inverse transfer functions representing the estimated social signal attenuation caused by each of the corresponding interagent pose parameters. The order in which the functions are applied is important for accurately estimating the speech output level.

First, the inverse attenuation function $A_{\alpha}^{-1}(L)$ of robot-to-human orientation (α) given a sound pressure level L must be applied. $A_{\alpha}^{-1}(L)$ is based on a transfer function representing how sound pressure levels measured by the robot's microphones vary with orientation to a sound source as a result of the robot platform and microphone configuration, similar to a human head-related transfer function (Gardner and Martin, 1994). This transfer function must be determined empirically for each robot platform using the process described in Gardner and Martin (1994), and is represented as a lookup table (there is no closed-form solution).

Next, the inverse attenuation function $A_d^{-1}(L')$ of interagent distance (d) given $L' = A_\alpha^{-1}(L)$ is applied. $A_d^{-1}(L')$ is calculated based on an equation for sound propagation:

$$A_d^{-1}(L') = L' + 20 \times \log(d/d_0) \quad (6.3)$$

where d_0 is the desired reference distance of SPL_{HR} .¹

Finally, the inverse attenuation function $A_\beta^{-1}(L'')$ of human-to-robot orientation (β) given $L'' = A_d^{-1}(L')$ is applied. $A_\beta^{-1}(L)$ is based on a “speaker directivity function” (Chu and Warnock, 2002) representing how sound pressure levels are attenuated by the body of the human speaker. There is no closed-form solution to human speaker directivity functions, so a well-established lookup table is used (Chu and Warnock, 2002).

The calculation of L'' using Equation 6.2 provides an estimate of *actual* SPL_{HR} at a measured *POS*. From this, the intercept (b), slope (m), and RSE (ϵ) of *predicted* SPL_{HR} in Equation 6.1 must be adapted. Given only a single value of *POS* (i.e., the first *POS* at which the robot and human ever interact), the only parameters that can be adapted are b (shifting the line up or down) and ϵ (widening or narrowing the estimated RSE). The robot must explore and interact at different interagent poses to adapt values of m , and the difference (ΔPOS) between these poses needs to be large enough to avoid over-fitting to a particular region of poses; appropriate values of ΔPOS will be investigated in the dissertation.

A graphical summary of the factors that influence the adaptation of SOL_{HR} over time is provided in Figure 6.1. The adaptation method would enable the robot to select interagent poses to improve its recognition of human speech based on more accurate estimates of SIL_{RH} .

¹ d_0 need only be specified once prior to robot deployment, and enables a normalized estimation of true SOL_{HR} (in dB SPL) at the mouth of any human speaker. In the dissertation, $d_0 = 0.1$ meters.

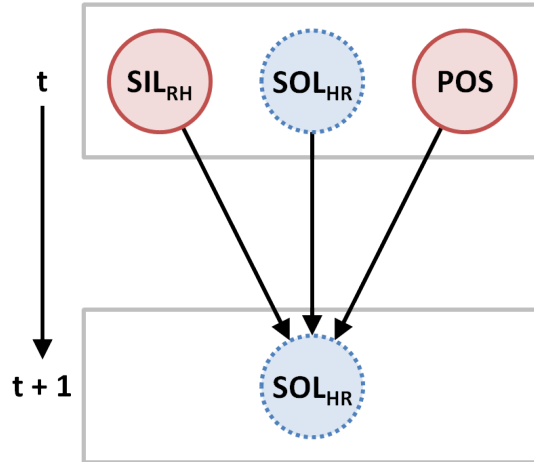


Figure 6.1: A graphical summary of the factors that influence the adaptation of human speech output levels (SOL_{HR}) over time. Red or blue indicates a state of either the robot or the human, respectively. A solid or dotted border indicates a state that is either measurable or latent (hidden), respectively.

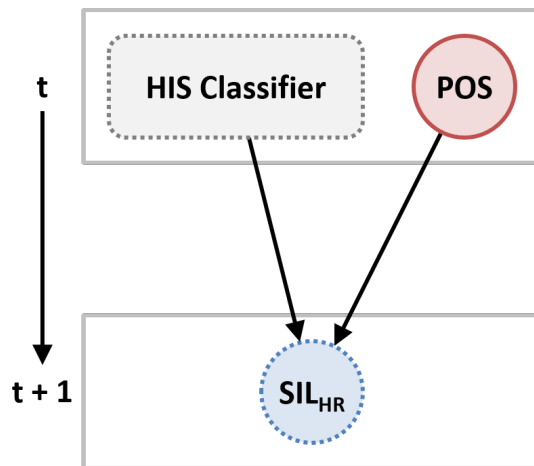


Figure 6.2: A graphical summary of the factors that influence the adaptation of human speech input levels (SIL_{HR}) over time. Red or blue indicates a state of either the robot or the human, respectively. The grey rounded box indicates a “hearing impairment/sensitivity” (HIS) classification system. A solid or dotted border indicates a state that is either measurable or latent (hidden), respectively.

6.2.2 Adaptive Models of Human Speech Input Levels (SIL_{HR})

In the implementation of the data-driven approach (Section 5.3), the robot selects parameters for its own speech output level (SOL_{RH}) that will likely maximize the ability of a human to hear and understand it (SIL_{HR}) (Section 5.3.1.1; Figure 5.4); however, neither the robot *nor the human* will perform as well if the models of human speech input are inaccurate. Thus, the objective of this extension is to adapt the models of human speech input levels SIL_{HR} .

The adaptation approach for human speech output levels (SOL_{HR}) described above (Section 6.2.1) relies on the ability of the robot to derive values of human speech *production* based on sensor (microphone) measurements; however, the robot has no direct mechanisms to detect how a human user is *perceiving* speech. Thus, human difficulty in hearing must be inferred from more complex social cues that go beyond the parameters represented in the data-driven approach². This would require the development of a “hearing impairment/sensitivity” (HIS) classifier, an automated recognition system for detecting human social signals that indicate hearing impairments or sensitivities.

A simple graphical representation of the proposed classifier and its influence on the model of SIL_{HR} over time is provided in Figure 6.2. The adaptation method would enable the robot to select speech output levels to improve human recognition of its speech based on more accurate estimates of SIL_{HR} .

²The necessary parameters would have to consider the *meaning* of the signal, while the data-driven approach focuses on the *manner* in which the signal is produced (Section 3.3.1).

6.2.3 Discussion

This section discussed challenges and proposed approaches for adapting human speech parameters in the unified framework of proxemics and multimodal communication established in this dissertation (Mead and Matarić, 2012, 2014). This extension of the framework would enable the robot to make better decisions about 1) its goal pose with respect to the user, improving robot speech and gesture recognition rates; and 2) its social signal production to the user, improving human speech understanding.

6.3 Conclusions

In the future, robots will be a part of our everyday lives. They will be in our workplaces and in our homes. They will be our collaborators and our companions. As they will be operating in human environments, it is crucial that these robots abide by natural social conventions to ensure their acceptance, adoption, and effectiveness. The work presented in this dissertation provides computational methods, models, and technologies that endow an autonomous robot with a subset of fundamental social capabilities that people often take for granted—where we position ourselves in face-to-face social encounters, how we adjust our speech and gestures to accommodate the needs of our interaction partners, and how we adapt our social behaviors in noisy and cluttered environments. This brings us one step closer to realizing a future in which humans and robots coexist.

Bibliography

- L. Adams and D. Zuckerman. The effect of lighting conditions on personal space requirements. *The Journal of General Psychology*, 118(4):335–340, 1991.
- J. Aiello. Human spatial behavior. In *Handbook of Environmental Psychology*, chapter 12. John Wiley & Sons, New York, New York, 1987.
- J. Aiello, D. Thompson, and D. Brodzinsky. How funny is crowding anyway? effects of group size, room size, and the introduction of humor. *Basic and Applied Social Psychology*, 4(2):192–207, 1981.
- M. Argyle and J. Dean. Eye-contact, distance, and affiliation. *Sociometry*, 28:289–304, 1965.
- J. Bailenson, J. Blascovich, A. Beall, and J. Loomis. Equilibrium theory revisited: Mutual gaze and personal space in virtual environments. *Presence*, 10(6):583–598, 2001.
- C. Bartneck, E. Croft, D. Kulic, and S. Zoghbi. Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots. *International Journal of Social Robotics*, 1(1):71–81, 2009.
- C. Breazeal. Social interactions in hri: The robot view. *IEEE Transactions on Man, Cybernetics and Systems*, 34(2):181–186, 2003.
- C. Breazeal. *Designing Sociable Robots*. MIT Press, Cambridge, Massachusetts, 2004.
- A. G. Brooks and R. C. Arkin. Behavioral overlays for non-verbal communication expression on a humanoid robot. *Autonomous Robots*, 22(1):55–74, 2007.
- J. Burgoon, L. Stern, and L. Dillman. *Interpersonal Adaptation: Dyadic Interaction Patterns*. Cambridge University Press, New York, New York, 1995.
- J. Cassell, J. Sullivan, and S. Prevost. *Embodied Conversational Agents*. MIT Press, Cambridge, Massachusetts, 2000.

- W. Chu and A. Warnock. *Detailed Directivity of Sound Fields Around Human Talkers*. Research Report (Institute for Research in Construction (Canada)). Institute for Research in Construction, 2002.
- K. Dautenhahn. Socially intelligent robots: Dimensions of human-robot interaction. *Philosophical Transactions of the Royal Society B*, 362:679–704, 2007.
- K. Dautenhahn, C. Nehaniv, M. Walters, B. Robins, H. Kose-Bagci, N. Mirza, and M. Blow. Kaspar a minimally expressive humanoid robot for human-robot interaction research. *Applied Bionics and Biomechanics*, 6(3–4):369–397, 2009.
- A. Dempster, N. Laird, and D. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of Royal Statistical Society*, 39:1–38, 1977.
- R. D. Deutsch. *Spatial Structurings in Everyday Face-to-face Behavior: A Neurocybernetic Model*. The Association for the Study of Man-Environment Relations, Orangeburg, 1977.
- G. Eresha, M. Haring, B. Endrass, E. Andre, and M. Obaid. Investigating the influence of culture on proxemic behaviors for humanoid robots. In *22nd IEEE International Symposium on Robot and Human Interactive Communication*, RO-MAN 2013, pages 430–435, 2013.
- G. Evans and R. Wener. Crowding and personal space invasion on the train: Please dont make me sit in the middle. *Journal of Environmental Psychology*, 27:90–94, 2007.
- J. Fasola and M. Matarić. Robot exercise instructor: A socially assistive robot system to monitor and encourage physical exercise for the elderly. In *19th IEEE International Symposium on Robot and Human Interactive Communication*, RO-MAN 2010, pages 416–421, 2010.
- J. Fasola and M. Matarić. Comparing physical and virtual embodiment in a socially assistive robot exercise coach for the elderly. Technical Report CRES-11-003, USC Center For Robotics and Embedded Systems, Los Angeles, California, 2011.
- D. Feil-Seifer and M. Matarić. Defining socially assistive robotics. In *International Conference on Rehabilitation Robotics*, ICRR’05, pages 465–468, Chicago, Illinois, 2005.
- D. Feil-Seifer and M. Matarić. B3ia: A control architecture for autonomous robot-assisted behavior intervention for children with autism spectrum disorders. In *The 17th IEEE International Symposium on Robot and Human Interactive Communication*, RO-MAN 2008, pages 328–333, 2008.

- D. Feil-Seifer and M. Matarić. Toward socially assistive robotics for augmenting interventions for children with autism spectrum disorders. *Experimental Robotics: Springer Tracts in Advanced Robotics*, 54:201–210, 2009.
- 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 *Proceedings of the 6th ACM/IEEE International Conference on Human-Robot Interaction*, HRI’11, pages 323–330, Lausanne, Switzerland, 2011a.
- D. Feil-Seifer and M. Matarić. People-aware navigation for goal-oriented behavior involving a human partner. In *IEEE International Conference on Development and Learning*, volume 2 of *ICDL’11*, pages 1–6, Frankfurt Am Main, Germany, 2011b.
- D. Fox, W. Burgard, and S. Thrun. The dynamic window approach to collision avoidance. *IEEE Robotics & Automation Magazine*, 4(1):23–33, 1997.
- B. Gardner and K. Martin. Hrtf measurements of a kemar dummy-head microphone. Technical Report 280, MIT Media Lab Perceptual Computing, Boston, Massachusetts, 1994.
- E. Geden and A. Begeman. Personal space preferences of hospitalized adults. *Research in Nursing and Health*, 4:237–241, 1981.
- B. Gerkey and K. Konolige. Planning and control in unstructured terrain. In *ICRA Workshop on Path Planning on Costmaps*, ICRA’08, Pasadena, California, 2008.
- E. Hall. A system for notation of proxemic behavior. *American Anthropologist*, 65:1003–1026, 1963.
- E. T. Hall. *The Silent Language*. Doubleday Company, New York, New York, 1959.
- E. T. Hall. *The Hidden Dimension*. Doubleday Company, Chicago, Illinois, 1966.
- E. T. Hall. *Handbook for Proxemic Research*. American Anthropology Association, Washington, D.C., 1974.
- L. Hayduk and S. Mainprize. Personal space of the blind. *Social Psychology Quarterly*, 43(2):216–223, 1980.
- H. Hediger. Studies of the psychology and behaviour of captive animals in zoos and circuses. *Butterworths Scientific Publications*, 1955.
- Z. Henkel, R. Murphy, and C. Bethel. Towards a computational method of scaling a robots behavior via proxemics. In *7th ACM/IEEE International Conference on Human-Robot Interaction*, pages 145–146, Boston, Massachusetts, 2012.

- Y. Hiroi and A. Ito. Influence of the size factor of a mobile robot moving toward a human on subjective acceptable distance. *Mobile Robots - Current Trends*, pages 177–190, 2011.
- H. Hüttenrauch, K. Eklundh, A. Green, and E. Topp. Investigating spatial relationships in human-robot interaction. In *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS'06*, pages 5052–5059, 2006.
- S. Jones and J. Aiello. Proxemic behavior of black and white first-, third-, and fifth-grade children. *Journal of Personality and Social Psychology*, 25(1):21–27, 1973.
- I. Kastanis and M. Slater. Reinforcement learning utilizes proxemics: an avatar learns to manipulate the position of people in immersive virtual reality. *Transactions on Applied Perception*, 9(1):1–15, 2012.
- A. Kendon. *Conducting Interaction - Patterns of Behavior in Focused Encounters*. Cambridge University Press, New York, New York, 1990.
- D. Kennedy, J. Glscher, J. Tyszka, and R. Adolphs. Personal space regulation by the human amygdala. *Nature Neuroscience*, 12:1226–1227, 2009.
- C. Kidd and C. Breazeal. Designing for long-term human-robot interaction and application to weight loss. Technical report, Massachusetts Institute of Technology, Boston, Massachusetts, 2011.
- E. Kim, E. Newland, R. Paul, and B. Scassellati. Robotic tools for prosodic training for children with asd: A case study. In *International Meeting for Autism Research, IMFAR'2008*, London, England, 2008.
- R. Kirby, J. Forlizzi, and R. Simmons. Natural person-following behavior for social robots. In *International Conference on Human-Robot Interaction*, Arlington, Virginia, 2007.
- R. Kirby, R. Simmons, and J. Forlizzi. Companion: A constraint optimizing method for person-acceptable navigation. In *Robot and Human Interactive Communication*, pages 607–6012, Toyama, Japan, 2009.
- D. Koller and N. Friedman. *Probabilistic Graphical Models*. MIT Press, Cambridge, Massachusetts, 2009.
- H. Kozima, C. Nakagawa, and Y. Yasuda. Children-robot interaction: A pilot study in autism therapy. *Progress in Brain Research*, 164:385–400, 2007.
- 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, HRI '10*, pages 267–274. IEEE Press, 2010.

- H. Kuzuoka, Y. Suzuki, J. Yamashita, and K. Yamazaki. Reconfiguring spatial formation arrangement by robot body orientation. In *HRI*, Osaka, Japan, 2010.
- B. Lawson. *Sociofugal and Sociopetal Space, The Language of Space*. Architectural Press, Oxford, 2001.
- M. Lee and L. Takayama. Now, i have a body: Uses and social norms for mobile remote presence in the workplace. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI'11, pages 33–42, Vancouver, British Columbia, 2011.
- A. Libin and J. Cohen-Mansfield. Therapeutic robocat for nursing home residents with dementia: Preliminary inquiry. *American Journal of Alzheimer's Disease and Other Dementias*, 19(2):111–116, 2004.
- D. M. Lloyd. Dont stand so close to me: The effect of auditory input on interpersonal space. *Perception*, 38(4):617–620, 2009.
- S. Low and D. Lawrence-Ziga. *The Anthropology of Space and Place: Locating Culture*. Blackwell Publishing, Oxford, 2003.
- T. W. Mallenby. The personal space of hard-of-hearing children after extended contact with “normals”. *British Journal of Social and Clinical Psychology*, 14(3):253–257, 1975.
- E. Marder-Eppstein, E. Berger, T. Foote, B. Gerkey, and K. Konolige. The office marathon: Robust navigation in an indoor office environment. In *IEEE International Conference on Robotics and Automation*, ICRA'10, pages 300–307, Anchorage, Alaska, 2010.
- N. Marquardt and S. Greenberg. Informing the design of proxemic interactions. *IEEE Pervasive Computing*, 11(2):14–23, 2012.
- M. Matarić, J. Eriksson, D. Feil-Seifer, and C. Winstein. Socially assistive robotics for post-stroke rehabilitation. *International Journal of NeuroEngineering and Rehabilitation*, 4(5):1–9, 2007.
- D. McNeill. *Hand and Mind: What Gestures Reveal about Thought*. Chicago University Press, 1992.
- D. McNeill. Gesture, gaze, and ground. In *Lecture Notes in Computer Science*. Springer-Verlag, 2005.
- R. Mead and M. J. Matarić. A probabilistic framework for autonomous proxemic control in situated and mobile human-robot interaction. In *7th ACM/IEEE International Conference on Human-Robot Interaction*, HRI'12, pages 193–194, Boston, Massachusetts, 2012.

- R. Mead and M. J. Matarić. Perceptual models of human-robot proxemics. In *14th International Symposium on Experimental Robotics, ISER'14*, page to appear, Marrakech/Essaouira, Morocco, 2014.
- R. Mead and M. J. Matarić. Robots have needs too: People adapt their proxemic behavior to improve autonomous robot recognition of human social signals. In *4th International Symposium on New Frontiers in Human-Robot Interaction, NF-HRI'15*, page to appear, Canterbury, United Kingdom, 2015.
- R. Mead, E. Wade, P. Johnson, A. S. Clair, S. Chen, and M. J. Mataric. An architecture for rehabilitation task practice in socially assistive human-robot interaction. In *Robot and Human Interactive Communication*, pages 404–409, 2010.
- 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, NordiCHI'12*, Copenhagen, Denmark, 2012.
- R. Mead, A. Atrash, and M. J. Matarić. Automated proxemic feature extraction and behavior recognition: Applications in human-robot interaction. *International Journal of Social Robotics*, 5(3):367–378, 2013.
- A. Mehrabian. *Nonverbal Communication*. Aldine Transcation, Piscataway, 1972.
- L. Morency, J. Whitehill, and J. Movellan. Generalized adaptive view-based appearance model: Integrated framework for monocular head pose estimation. In *8th IEEE International Conference on Automatic Face Gesture Recognition, FG'08*, pages 1–8, 2008.
- J. Mumm and B. Mutlu. Human-robot proxemics: Physical and psychological distancing in human-robot interaction. In *6th ACM/IEEE International Conference on Human-Robot Interaction, HRI-2011*, pages 331–338, Lausanne, 2011.
- A. Murata. Basic characteristics of human's distance estimation. In *1999 IEEE International Conference on Systems, Man, and Cybernetics*, volume 2 of *SMC'99*, pages 38–43, 1999.
- J. Nelder and R. Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society*, 135(3):370–384, 1972.
- T. Oosterhout and A. Visser. A visual method for robot proxemics measurements. In *HRI Workshop on Metrics for Human-Robot Interaction*, Amsterdam, 2008.
- A. Ozyurek. 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):688–704, 2002.

- K. Pearsons, R. Bennett, and S. Fidell. *Speech Levels in Various Noise Environments*. U.S. Environmental Protection Agency, Washington, District of Columbia, 1977.
- G. Price and J. Dabbs Jr. Sex, setting, and personal space: Changes as children grow older. *Personal Social Psychology Bulletin*, 1:362–363, 1974.
- M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Ng. Ros: An open-source robot operating system. In *ICRA Workshop on Open Source Software*, ICRA’09, Kobe, Japan, 2009.
- L. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. In *Readings in Speech Recognition*, pages 267–296, 1990.
- I. Rae, L. Takayama, and B. Mutlu. The influence of height in robot-mediated communication. In *8th ACM/IEEE International Conference on Human-Robot Interaction, HRI-2013*, pages 1–8, Tokyo, Japan, 2005.
- D. Ricks and M. Colton. Trends and considerations in robot-assisted autism therapy. In *IEEE International Conference on Robotics and Automation, ICRA’10*, pages 4354–4359, Anchorage, Alaska, 2010.
- N. Rossini. The analysis of gesture: establishing a set of parameters. In *Gesture-based Communication in Human-Computer Interaction*, volume 2915, pages 463–464. Springer-Verlag, 2004.
- 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*, pages 109–116, 2009.
- B. Scassellati. How social robots will help us to diagnose, treat, and understand autism. *Robotics Research, Springer Tracts in Advanced Robotics*, 28:552–563, 2007.
- E. Schegloff. Body torque. *Social Research*, 65(3):535–596, 1998.
- H. Schöne. *Spatial Orientation: The Spatial Control of Behavior in Animals and Man*. Princeton University Press, Princeton, 1984.
- C. Shi, M. Shimada, T. Kanada, H. Ishiguro, and N. Hagita. Spatial formation model for initiating conversation. In *RSS*, Los Angeles, 2011.
- E. Short, J. Hart, M. Vu, and B. Scassellati. No fair!! an interaction with a cheating robot. In *5th ACM/IEEE International Conference on Human-Robot Interaction*, pages 219–226, Osaka, Japan, 2010.
- J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. In *CVPR*, 2011.

- R. Sommer. Sociofugal space. *The American Journal of Sociology*, 72(6):654–660, 1967.
- S. Stevens. On the psychological law. *Psychological Review*, 64:153–181, 2007.
- L. Takayama and C. Pantofaru. Influences on proxemic behaviors in human-robot interaction. In *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, IROS’09, pages 5495–5502, 2009.
- A. Tapus, M. Matarić, and B. Scassellati. The grand challenges in socially assistive robotics. *IEEE Robotics and Automation Magazine*, 14(1):35–42, 2007.
- A. Tapus, C. Tapus, and M. Matarić. The use of socially assistive robots in the design of intelligent cognitive therapies for people with dementia. In *Proceedings of the International Conference on Rehabilitation Robotics*, ICORR’09, pages 924–929, 2009a.
- A. Tapus, C. Tapus, and M. Matarić. Music therapist robot for people suffering from dementia: Longitudinal study. In *Proceedings of the International Conference on Alzheimer’s Disease*, ICAD’09, pages 266–266, 2009b.
- 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*, ICSR’11, pages 21–30, 2011.
- E. Torta, R. H. Cuijpers, and J. F. Juola. Design of a parametric model of personal space for robotic social navigation. *International Journal of Social Robotics*, 5(3): 357–365, 2013.
- H. Traunmüller and A. Eriksson. Acoustic effects of variation in vocal effort by men, women, and children. *Journal of the Acoustical Society of America*, 107(6):3438–3451, 2000.
- P. Trautman and A. Krause. Unfreezing the robot: Navigation in dense, interacting crowds. In *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, IROS’10, pages 797–803, 2010.
- K. Wada and T. Shibata. Living with seal robots—its sociopsychological and physiological influences on the elderly at a care house. *IEEE Transactions on Robotics*, 23(5):972–980, 2007.
- K. Wada, T. Shibata, T. Saito, and K. Tanie. Effects of robot-assisted activity for elderly people and nurses at a day service center. *Proceedings of the IEEE*, 92(11): 1780–1788, 2004.
- J. Wainer, E. Ferrari, K. Dautenhahn, and B. Robins. The effectiveness of using a robotics class to foster collaboration among groups of children with autism in an exploratory study. *Personal and Ubiquitous Computing*, 14(5):445–455, 2010.

- M. Walters, D. Syrdal, K. Koay, K. Dautenhahn, and R. te Boekhorst. Human approach distances to a mechanical-looking robot with different robot voice styles. In *The 17th IEEE International Symposium on Robot and Human Interactive Communication*, RO-MAN 2008, pages 707–712, 2008.
- 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*, pages 144–149, Edinburgh, 2009.
- J. D. Webb and M. J. Weber. Influence of sensory abilities on the interpersonal distance of the elderly. *Environment & Behavior*, 35(5):695–711, 2003.
- W. Welch, U. Lahiri, Z. Warren, and N. Sarkar. An approach to the design of socially acceptable robots for children with autism spectrum disorders. *International Journal of Social Robotics*, 2:391–403, 2010.
- M. Wöelfel and J. McDonough. *Distant Speech Recognition*. John Wiley & Sons Ltd, West Sussex, United Kingdom, 2009.