

Toward Robot Adaptation of Human Speech and Gesture Parameters in a Unified Framework of Proxemics and Multimodal Communication

Ross Mead¹ and Maja J Matarić²

Abstract—In this work, we present our unified framework of proxemics and multimodal communication based on static data-driven models of how social signals (speech and gesture) are produced and perceived by both a social robot and a human user, and propose extensions to enable dynamic model adaptation of these human communication parameters.

I. INTRODUCTION

During face-to-face human-robot interaction (HRI), a *social robot* 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 most signals, these *social signals* are attenuated by distance; this influences how the signals are *produced* (e.g., humans adapt to increased distance by talking louder [1]) and *perceived* (e.g., sensory sensitivities [2] or impairments [3], [4] in socially assistive domains [5]) by both the human and the robot. Thus, for a *mobile* social robot, proper selection of interagent pose (proxemics) and multimodal communication parameters can be a deciding factor in successful HRI [6].

Proxemics is the study of the interpretation, manipulation, and dynamics of human spatial behavior in co-present social encounters [1]. Proxemic behavior is governed by sociocultural norms, which, in effect, determine the overall sensory experience of each interacting participant [1]. The distance and orientation between two social agents influences how social stimuli (e.g., speech and gestures) are produced and perceived [1], [6]; changes in these spatial parameters offer cues that can indicate difficulties in understanding social stimuli. Autonomously recognizing these often subtle and noisy cues can provide a means of adapting to the needs of a human user with regard to multimodal communication in face-to-face HRI [7].

In this work, we present our unified framework of proxemics and multimodal communication [6], [8] based on static data-driven models of how social signals (speech and gesture) are produced and perceived by both the human and the robot, and propose extensions to enable dynamic model adaptation of these human communication parameters.

*This work is supported in part by an NSF Graduate Research Fellowship and the NSF National Robotics Initiative (IIS-1208500).

¹Ross Mead is a PhD Candidate in Computer Science, University of Southern California, Los Angeles, CA 90089, USA rossmead@usc.edu

²Maja J Matarić is a Professor of Computer Science, Neuroscience, and Pediatrics, University of Southern California, Los Angeles, CA 90089, USA mataric@usc.edu

II. BACKGROUND

In developing an approach for adaptation of human speech and gesture, we first consider related work in modeling proxemics in both human-human and human-robot interactions, and the relationship between proxemics and multimodal communication; we then review user adaptation in HRI.

A. Model Representations of Proxemics

There exists a rich body of work in the social sciences that seeks to explain proxemic phenomena. Some of these proxemic models have been investigated in human-robot interaction scenarios. In [9], we grouped these approaches into three categories of representation: 1) *physical*, 2) *psychological*, and 3) *psychophysical*. Physical and psychological representations have received the most attention in both human-human and human-robot interactions; however, these representations often struggle to model the complex dynamics between interacting agents and environmental interference (e.g., loud noise or visual occlusions). In [9], we proposed the use of a psychophysical representation of proxemic behavior for modeling the situatedness of the agents within the environment, and to serve as a bridge between the physical and psychological proxemic representations.

1) *Physical Representations of Proxemics*: *Physical* representations of proxemics are concerned with how space is occupied by two or more bodies, relating these bodies via low-level spatial parameters of distance and orientation [10], [11] (Figure 1). These representations are the most common in the analysis of proxemics in both human-human and human-machine interactions. Virtual humans have been shown to manipulate the position of people in immersive virtual environments [12]. Rule-based spatial interaction controllers have also been applied to human-robot interaction (HRI) [13], [14], [15]. Marquardt & Greenberg [16] expanded upon physical proxemic representation for ubiquitous computing by adding movement, identity, and location.

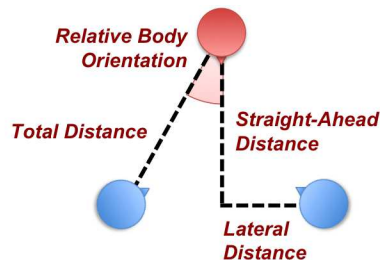


Fig. 1. The *physical* proxemic representation, such as in [10], is the most common in the analysis of proxemic behavior.

2) *Psychological Representations of Proxemics*: Psychological representations of proxemics are concerned with the high-level interpersonal relationship between two or more agents. For example, Affiliative Conflict Theory [17] proposes that interagent spacing is governed by an equilibrium of intimacy, consisting of the amount of mutual eye gaze and smiling, the intimacy of topic or thought, and posture and arm configurations. Interpersonal Adaptation Theory [18] extends this by considering more dynamic behaviors—such as the amount of convergence, divergence, mirroring, compensation, and reciprocity—between social partners. These models have been investigated in HRI [19], [20].

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 pose (Figure 2). This dynamic can be expressed using a psychophysical representation, situating the agents in the environment [8], [9].

3) *Psychophysical Representations of Proxemics*: Psychophysical representations are concerned with the perception and production of social stimuli by two or more interacting agents. Edward T. Hall [21] coined the term proxemics, and proposed culture-defined zones of intimate, personal, social, and public space (a psychological representation) [1]; however, Hall [1], [22] posited that these zones are characterized by psychophysical factors, such as the visual (field of view), auditory (voice loudness), olfactory (smell), thermal (body heat), and somatosensory (tactile and kinesthetic) experiences of each interacting participant. The effect of sensory experience on proxemic behavior is evident in the analysis of the personal space requirements of people with hearing [3] and visual [4] impairments.

Psychophysical representations relate sensory experience of multimodal communication (social stimuli, such as speech and gesture) to spatial parameters (such as distance and orientation); however, these representations have yet to be thoroughly examined, formalized, and adopted in HRI.

In previous work [7], we presented a method for both physical and psychophysical proxemic feature extraction inspired by literature, and demonstrated increased performance of the psychophysical representation over the physical representation in the automated recognition of human spatial dynamics for initiating and terminating social interactions. In [6], we utilized this psychophysical philosophy to model social robot performance attenuated by distance, which was used to implement an autonomous robot controller for proxemic behavior and multimodal that attempts to maximize the perception of speech and gesture for both the robot and the human in face-to-face HRI. In [23], we demonstrated that users naturally adapted their proxemic behavior to compensate for poor robot understanding of their social signals in at a variety of interaction conditions. We now consider methods in which a social robot can dynamically adapt its own behaviors to accommodate for the multimodal communication needs of a human user.

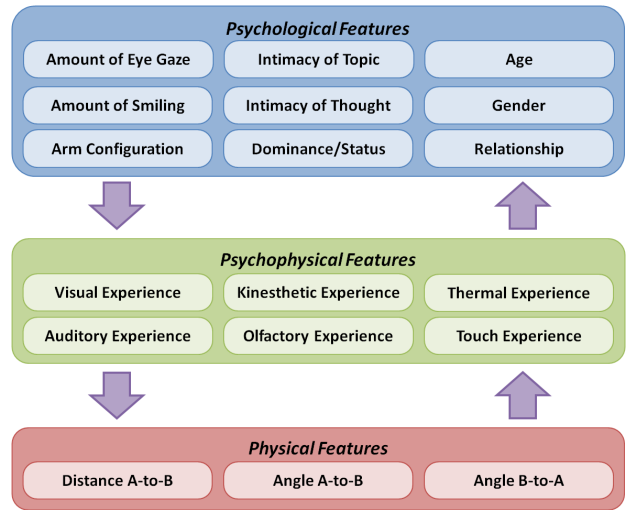


Fig. 2. Psychological factors dictate the desired psychophysical (sensory) experience of each agent, which is manifested physically through the manipulation of space via change in position and orientation [7], [9].

B. User Adaptation in Human-Robot Interaction

We focus our review of work in adaptation with regard to our interaction parameters of interest (pose, speech, and gesture), particularly in socially assistive domains [5].

Mitsunaga *et al.* [24] utilized policy gradient reinforcement learning to adapt user preferences of proxemic (spacing), kinesic (body movement), oculosic (eye gaze), chronemic (timing) behaviors [25]. Hasanuzzaman & Ueno [26] uniquely identified user faces, then used a multi-clustering approach to learn user-specific mappings from static human face, hand, and arm gestures to robot behaviors. These bodies of work focus on the *desires* or preferences of human users; our work distinguishes itself in its focus on user *needs*.

User personalization is particularly important in the field of *socially assistive robotics* (SAR), in which robots interact with people with special needs to improve quality of life [5]. François *et al.* [27] demonstrated an approach to online adaptation of interaction styles (“strong” vs. “gentle” play) with children with autism spectrum disorders (ASDs) using self-organizing maps; similarly, Feil-Seifer & Matarić [2] used Gaussian mixture models to classify positive and negative reactions of children with ASDs during an interaction with a SAR system, which could be used to adapt robot behavior for the child (proposed as future work). Fasola & Matarić [28] performed a longitudinal study with a SAR exercise coach that personalized its training regiment with older adults based on its history of interactions with them. We highlight these two populations as they are often characterized by sensory sensitivities and impairments, respectively, which are of particular interest in our own adaptation work.

In Section III, we present our unified model (based on a psychophysical representation [8], [9]; discussed in Section II-A.3) of proxemics and multimodal communication [6]; however, this model is currently static. In Section IV, we discuss challenges and propose approaches for dynamically adapting speech and gesture parameters of a human user.

III. UNIFIED FRAMEWORK OF PROXEMICS AND MULTIMODAL COMMUNICATION

Our probabilistic proxemic framework considers how all represented social agents—both humans and robots—experience a co-present interaction [8]. We model the production (output) and perception (input) of speech and gesture conditioned on interagent pose (position and orientation).

A. Parameter Definitions

Consider two social agents, A and B , that are co-located and would like to interact.¹ At any point in time and from any location in the environment, an agent must be capable of estimating the following interaction parameters:

- 1) an **interagent pose**, POS —At what position and orientations will B stand relative to A ?
- 2) a **speech output level**, SOL_{BA} —How loudly will B speak for A to hear?
- 3) a **gesture output level**, GOL_{BA} —In what volume of space will B gesture for A to see?
- 4) a **speech input level**, SIL_{AB} —How will A hear speech produced by B ?
- 5) a **gesture input level**, GIL_{AB} —How will A see gestures produced by B ?

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

Interagent pose (POS) is expressed as a face-to-face distance (d) and two orientations—one from A to B (α), and one from B to A (β). These pose parameters attenuate speech and gesture output produced by agents, which, in turn, impacts perceived speech and gesture input levels. Note that POS of the robot R cannot be changed in isolation, as it is relative to the often changing pose of the human H , detected by a person-tracking system (e.g., a Microsoft Kinect).

Speech output and input levels (SOL and SIL) are each represented as a sound pressure level (in dB SPL), a logarithmic measure of sound pressure relative to ambient noise in the environment, serving as a signal-to-noise ratio. SOL is generated by human speech and robot speakers. SIL is detected by human ears and robot microphones, and is further attenuated by human equal-loudness contours [29] and the robot speech recognition system, respectively.

Gesture output and input levels (GOL and GIL) are each represented as a 3D volume of space gestures will occupy, referred to as a *gesture locus* [30], [31]. GOL is the locations of body parts (we focus on the face and hands) in physical space during gesture production, and is based on the physical embodiment and kinematics of both the human and the robot [32], which can vary with distance [33] and orientation [34] between agents. GIL is the region of the human retinal field (foveal, macular, or peripheral [22]) or the robot gesture recognition system (e.g., using a Microsoft Kinect) occupied by GOL , and is attenuated based on human visual acuity [22], [4] and robot camera resolution.

¹Our current work focuses on dyadic interactions (i.e., between two agents) using only speech and gesture; however, our principled framework [1], [22] is extensible to any number of social agents and social signals [8].

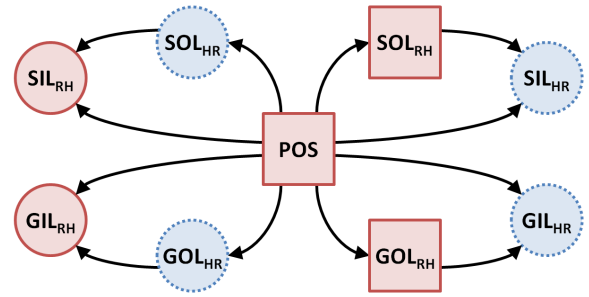


Fig. 3. A Bayesian network topology relating pose, speech, and gesture parameters. Red or blue indicates a parameter of either the robot or the human, respectively. A circle or square indicates either a state or an action, respectively. A solid or dotted border indicates a quantity that is either measurable/controllable or latent (hidden), respectively.

B. Parameter Modeling

We model distributions over these pose, speech, and gesture parameters and their relationships as a Bayesian network to represent: 1) how an agent will pose itself relative to the other, 2) how interagent pose influences agent speech and gesture output (production), and 3) how interagent pose influences agent speech and gesture input (perception); we denote the robot as R and the human as H (Figure 3).

In [6], we performed data collections (in both human-human and human-robot interactions) to inform the parameters of the proposed Bayesian network. We implemented the resulting models in an autonomous control system for a social robot [6], [8]; the controller enables the robot to 1) select an interagent pose that *attempts* to maximize the robot’s understanding of human social signals, and 2) select robot speech and gesture output levels that *attempt* to maximize human understanding of its social signals. Our current models and controller are based on an averaging of results *across* participants; however, the models do not account for variability *between* participants. Implications of these two approaches are discussed below.

Across participants (i.e., considered as a collective), we identified a strong positive linear relationship between human speech output levels and distance to both human and robot interaction partners (though people spoke louder to the robot at near distances than with another person). We observed no statistically significant differences in human gesture output levels at different distances (though this is believed to have been due to the simple task they were performing).

Between participants (i.e., considered as individuals), we noticed significant variability 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; our models of robot speech and gesture output are currently static and do not accommodate for these differences.

These insights led us to consider methods for dynamically adapting our models of human speech and gesture production and perception attenuated by interaction pose.

IV. TOWARD ADAPTATION OF HUMAN SPEECH AND GESTURE PARAMETERS

The limitations of our static model approach highlighted in Section III-B led us to ask: how can the robot dynamically adapt its models of 1) *human speech* and *gesture output levels*; and 2) *human speech* and *gesture input levels*? These questions, as well as each communication channel (speech and gesture) within each question, are characterized by distinct implications and challenges—subsequently, each warrants its own unique solution. We propose approaches to the adaptation these human communication output (Section IV-A) and input (Section IV-B) parameters below, which we are implementing and testing in our ongoing work.

A. Adaptation of Human Speech and Gesture Output Levels

In our framework [6], [8], the robot selects and moves to a pose that will most likely maximize its perceptions of social signals (SIL_{RH} and GIL_{RH}) based on models of human speech and gesture output (SOL_{HR} and GOL_{HR}) (Figure 3); thus, the robot will not perform optimally if the models of human multimodal output are inaccurate. Our goal then is to dynamically adapt these models of SOL_{HR} and GOL_{HR} ; a graphical summary of the factors that influence these models and their adaptation over time is provided in Figure 4.

1) *Human Speech Output Level (SOL_{HR})*: As noted in Section III-B, human speech output levels (SOL_{HR}) can be estimated by a linear function of interagent pose (POS) [6]:

$$SOL_{HR}(POS) = b + m \times POS + \epsilon \quad (1)$$

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

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

$$A^{-1}(SIL_{RH}) = A_{\beta}^{-1}(A_d^{-1}(A_{\alpha}^{-1}(SIL_{RH}))) \quad (2)$$

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, especially for adaptation of speech pitch and energy [35] (future work).

First, we apply the inverse attenuation function $A_{\alpha}^{-1}(L)$ of robot-to-human orientation (α) given a sound pressure level L . $A_{\alpha}^{-1}(L)$ is based on a transfer function representing how sound pressure levels measured by the robot’s microphones will 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 [36]). This transfer function must be determined empirically for each robot platform (using the process in [36]) and is represented as a lookup table (there is no closed-form solution).

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

$$A_d^{-1}(L') = L' + 20 \times \log(d/d_0) \quad (3)$$

where d_0 is the desired reference distance of SPL_{HR} .²

Last, we apply the inverse attenuation function $A_{\beta}^{-1}(L'')$ of human-to-robot orientation (β) given $L'' = A_d^{-1}(L')$. $A_{\beta}^{-1}(L)$ is based on a “speaker directivity function” [37] 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 (available in [37]).

The calculation of L'' using Equation 2 gives us an estimate of *actual* SPL_{HR} at a measured POS . From this, we need to adapt the intercept (b), slope (m), and RSE (ϵ) of *predicted* SPL_{HR} in Equation 1. Given only a single value of POS (i.e., the first POS at which the robot and human ever interact), we can only adapt b (shifting the line up or down) and ϵ (widening or narrowing the estimated RSE). The robot needs to 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 overfitting to a particular region of poses; we are investigating appropriate values of ΔPOS in our ongoing work.

2) *Human Gesture Output Level (GOL_{HR})*: As noted in Section III-B, we identified no statistically significant relationship between human gesture output level (GOL_{HR}) and interagent pose (POS) [6]. Using person-tracking technologies (e.g., the Microsoft Kinect), we can accurately detect the locations of body parts used in human gesture production. While gesture recognition performance is attenuated by distance due to camera resolution, the attenuation effect on the actual gesture output is negligible—unlike speech, an attenuation function is not necessary. This simplifies the adaptation of GOL_{HR} based on GIL_{RH} —during human gesturing, we continuously update distributions $\mathcal{N}(\mu, \sigma)$ of the minimum and maximum locations of detected body parts (the face and hands) along each axis (relative to the user).

The proposed methods of adapting SOL_{HR} and GOL_{HR} would enable the robot to select poses that improve its recognition of human speech and gesture based on more accurate estimates of SIL_{RH} and GIL_{RH} , respectively.

² d_0 need only ever 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 our case, $d_0 = 0.1$ meters.

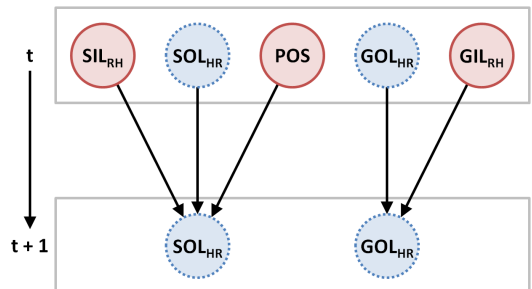


Fig. 4. A graphical summary of the factors that influence our adaptation of human speech and gesture output levels over time. Red or blue indicates a state of either the robot or the human, respectively. A solid or dotted border indicates a state that is either measurable or latent (hidden), respectively.

B. Adaptation of Human Speech and Gesture Input Levels

In our framework [6], the robot selects parameters for its own speech and gesture output levels (SOL_{RH} and GOL_{RH}) that will likely maximize human perceptions of its social signals (SIL_{HR} and GIL_{HR}) (Figure 3); however, neither the robot *nor the human* will perform optimally if the models of human multimodal input are inaccurate. Our goal then is to adapt these models of SIL_{HR} and GIL_{HR} .

Our adaptation approaches for human speech and gesture output levels (Section IV-A) rely on the ability of the robot to derive adapted values of social signals based on sensor measurements; however, we have no direct mechanisms to detect how a human user is perceiving speech and gesture—we must infer it from more complex social cues that go beyond the parameters currently represented in the model. Thus, we must now consider the development of automated recognition systems for independently classifying social signals that indicate 1) hearing impairments or sensitivities (HIS), or 2) visual impairments or sensitivities (VIS); these classifiers can be used by the robot to adapt its models of human speech and gesture input (SIL_{HR} and GIL_{HR}).

For each classifier, we must consider social signals that are 1) *communicative* (produced with intent) vs. *informative* (produced without intent), and 2) *state-based* (describing the impairment or sensitivity) vs. *action-based* (describing the desired robot action to resolve the impairment or sensitivity). In our work, we are primarily concerned with spatial, aural (speech), and visual (gesture) cues within these categories, though other social cues could be considered. A simple graphical representation of the proposed classifiers and their influence on the models over time is provided in Figure 5; examples (not a complete list) of social cues that could be used by each classifier (HIS and VIS) are proposed below.

1) *Human Speech Input Level (SIL_{HR})*: A variety of social cues could indicate a hearing impairment or sensitivity, which would need to be classified by the classifier (HIS) to adapt human speech input levels (SIL_{HR}). For example, a human user could say, “I cannot hear you.” (communicative; state-based) or “Please speak up.” (communicative; action-based), indicating difficulty hearing. A turn of the head or cupping of the hand around the ear (informative or communicative; state-based) orients the user’s ear to the speaker; a user who moves closer to the robot (informative; state-based) might also suffer from a hearing impairment, as shown with children who are deaf [4] and older adults [38].

Sound sensitivities are rather common with children with ASDs, but can also be relevant in particular user interaction contexts (e.g., a robot in a library). Loud or sudden noises from a social robot can trigger negative user responses, as was shown with children with ASDs who distanced themselves from the robot when startled (informative; state-based) [2]. Covering of the ears (informative or communicative; state-based), a finger held in front of the mouth (communicative; action-based), or a simple shushing utterance (communicative; action-based) can be classified by HIS to indicate a need to dampen the models of SIL_{HR} .

2) *Human Gesture Input Level (GIL_{HR})*: Social cues indicating visual impairments or sensitivities are often less salient and less clear than that of hearing discrepancies, which will increase uncertainty in the classifier VIS . For example, a human user might move closer to observe details [4], might move farther away to gain a broader perspective [38], or look or face in atypical directions to compensate for distorted vision or macular degeneration [38] (informative or communicative; state-based)—these spatial cues appear to conflict with each other. More clear signals of visual impairments include squinting of the eyes (informative; state-based) or simply saying, “I cannot see.” (communicative; state-based). A robust visual impairment/sensitivity (VIS) classifier would need to accommodate each of these factors in predicting changes in human gesture input levels GIL_{HR} .

Fast or sudden movements by the robot could trigger negative responses in humans with visual sensitivities; both children with ASDs [2] and older adults (especially those with limited mobility) [38] will increase their distance from a social partner (informative; state-based) moving quickly or broadly. A covering of the eyes (informative or communicative; state-based), a beat gesture with palms down (informative or communicative; action-based) [31], or a “calm down” utterance (communicative; action-based) could be classified by the VIS classifier as overwhelming visual stimuli, necessitating a dampening of GIL_{HR} , and, subsequently, a reduction in robot gesture output (GOL_{RH}).

The proposed methods of adapting SIL_{HR} and GIL_{HR} would enable the robot to select multimodal communication parameters (SOL_{RH} and GOL_{RH}) that improve human perception of the robot’s social signals.

The provided examples of social cues that could be used by each classifier (HIS and VIS) are a subset of a much more broad and complex space, which we intend to explore and model in data collections with human users. We are planning a large-scale data collection with visitors of an open house event at the University of Southern California. Participants will interact with a social robot in a scenario that will explore the space of social cues indicating hearing and visual impairments and sensitivities. From these interactions, we will identify and model distributions over categories of social cues relevant to our classifiers (HIS and VIS).

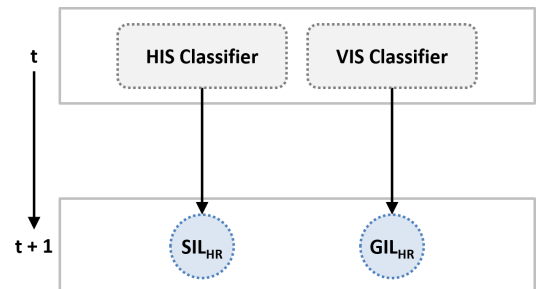


Fig. 5. A graphical summary of the factors that influence our adaptation of human speech and gesture input levels over time. A blue circle indicates a human state. A grey rounded box indicates an automated recognition system (i.e., a classifier). A dotted border indicates a state that is latent (hidden).

V. SUMMARY AND FUTURE WORK

We have discussed challenges and proposed approaches for adapting human speech and gesture parameters in our unified framework of proxemics and multimodal communication [6], [8]. The extension of our 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 and gesture understanding. We are integrating the proposed adaptation techniques into our framework, and will experimentally explore, inform, and validate each approach in our ongoing and future work.

This research establishes a foundational component of HRI, enabling the development of robust decision-making and control systems for social robots. Furthermore, this work has implications for technology personalization and will be evaluated in socially assistive contexts with people with hearing or visual impairments or sensitivities.

REFERENCES

- [1] E. T. Hall, *The Hidden Dimension*. Chicago, Illinois: Doubleday Company, 1966.
- [2] 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*, ser. HRI'11, Lausanne, Switzerland, 2011, pp. 323–330.
- [3] T. W. Mullenby, "The personal space of hard-of-hearing children after extended contact with "normals"," *British Journal of Social and Clinical Psychology*, vol. 14, no. 3, pp. 253–257, 1975.
- [4] L. Hayduk and S. Mainprize, "Personal space of the blind," *Social Psychology Quarterly*, vol. 43(2), pp. 216–223, 1980.
- [5] A. Tapus, M. Matarić, and B. Scassellati, "The grand challenges in socially assistive robotics," *IEEE Robotics and Automation Magazine*, vol. 14, no. 1, pp. 35–42, 2007.
- [6] R. Mead and M. J. Matarić, "Perceptual models of human-robot proxemics," in *14th International Symposium on Experimental Robotics*, ser. ISER'14, Marrakech/Essaouira, Morocco, 2014, p. to appear.
- [7] 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*, vol. 5, no. 3, pp. 367–378, 2013.
- [8] 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*, ser. HRI'12, Boston, Massachusetts, 2012, pp. 193–194.
- [9] 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*, ser. NordiCHI'12, Copenhagen, Denmark, 2012.
- [10] A. Mehrabian, *Nonverbal Communication*. Piscataway: Aldine Transaction, 1972.
- [11] E. Schegloff, "Body torque," *Social Research*, vol. 65, no. 3, pp. 535–596, 1998.
- [12] 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*, vol. 9, no. 1, pp. 1–15, 2012.
- [13] H. Kuzuoka, Y. Suzuki, J. Yamashita, and K. Yamazaki, "Reconfiguring spatial formation arrangement by robot body orientation," in *HRI, Osaka, Japan*, 2010.
- [14] M. Walters, K. Dautenhahn, R. Boekhorst, K. Koay, D. Syrdal, and C. Nehaniv, "An empirical framework for human-robot proxemics," in *New Frontiers in Human-Robot Interaction*, Edinburgh, 2009, pp. 144–149.
- [15] H. Hüttenrauch, K. Eklundh, A. Green, and E. Topp, "Investigating spatial relationships in human-robot interaction," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, ser. IROS'06, Beijing, China, 2006.
- [16] N. Marquardt and S. Greenberg, "Informing the design of proxemic interactions," *IEEE Pervasive Computing*, vol. 11, no. 2, pp. 14–23, 2012.
- [17] M. Argyle and J. Dean, "Eye-contact, distance, and affiliation," *Sociometry*, vol. 28, pp. 289–304, 1965.
- [18] J. Burgoon, L. Stern, and L. Dillman, *Interpersonal Adaptation: Dyadic Interaction Patterns*. New York, New York: Cambridge University Press, 1995.
- [19] 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*, ser. HRI-2011, Lausanne, 2011, pp. 331–338.
- [20] L. Takayama and C. Pantofaru, "Influences on proxemic behaviors in human-robot interaction," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, ser. IROS'09, 2009, pp. 5495–5502.
- [21] E. T. Hall, *The Silent Language*. New York, New York: Doubleday Company, 1959.
- [22] E. Hall, "A system for notation of proxemic behavior," *American Anthropologist*, vol. 65, pp. 1003–1026, 1963.
- [23] 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*, ser. NF-HRI'15, Canterbury, United Kingdom, 2015, p. to appear.
- [24] N. Mitsunaga, C. Smith, T. Kanda, H. Ishiguro, and N. Hagita, "Adapting robot behavior for human-robot interaction," *IEEE Transactions on Robotics*, vol. 24, no. 4, pp. 911–916, 2008.
- [25] R. Mead, A. Atrash, E. Kaszubski, A. St. Clair, J. Greczek, C. Clabaugh, B. Kohan, and M. J. Matarić, "Building blocks of social intelligence: Enabling autonomy for socially intelligent and assistive robots," in *Association for the Advancement of Artificial Intelligence Fall Symposium on Artificial Intelligence and Human-Robot Interaction*, ser. AI-HRI, Arlington, Virginia, 2014.
- [26] M. Hasanuzzaman and H. Ueno, "User, gesture and robot behaviour adaptation for human-robot interaction," *The Future of Humanoid Robots - Research and Applications*.
- [27] D. François, D. Polani, and K. Dautenhahn, "On-line behaviour classification and adaptation to human-robot interaction styles," in *2nd ACM/IEEE International Conference on Human-Robot Interaction*, ser. HRI'07, Arlington, Virginia, 2007, pp. 295–302.
- [28] J. Fasola and M. Matarić, "A socially assistive robot exercise coach for the elderly," *Journal of Human-Robot Interaction*, vol. 2, no. 2, pp. 3–32, 2013.
- [29] ISO, "Acoustics – normal equal-loudness-level contours (iso 226:2003)," *International Organization for Standardization*, 2003.
- [30] N. Rossini, "The analysis of gesture: establishing a set of parameters," in *Gesture-based Communication in Human-Computer Interaction*, vol. 2915. Springer-Verlag, 2004, pp. 463–464.
- [31] D. McNeill, *Gesture and Thought*. Chicago University Press, 2005.
- [32] R. Mead, E. Wade, P. Johnson, A. S. Clair, S. Chen, and M. J. Matarić, "An architecture for rehabilitation task practice in socially assistive human-robot interaction," in *Robot and Human Interactive Communication*, 2010, pp. 404–409.
- [33] A. G. Brooks and R. C. Arkin, "Behavioral overlays for non-verbal communication expression on a humanoid robot," *Autonomous Robots*, vol. 22, no. 1, pp. 55–74, 2007.
- [34] 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*, vol. 46, no. 4, pp. 688–704, 2002.
- [35] J. Rosenberg and J. Hirschberg, "Detecting pitch accents at the word, syllable and vowel level," in *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Short Papers*, 2009, pp. 81–84.
- [36] B. Gardner and K. Martin, "Hrtf measurements of a kemar dummy-head microphone," MIT Media Lab Perceptual Computing, Boston, Massachusetts, Tech. Rep. 280, 1994.
- [37] W. Chu and A. Warnock, *Detailed Directivity of Sound Fields Around Human Talkers*, ser. Research Report (Institute for Research in Construction (Canada)). Institute for Research in Construction, 2002.
- [38] J. D. Webb and M. J. Weber, "Influence of sensory abilities on the interpersonal distance of the elderly," in *Environment & Behavior*, vol. 35, no. 5, 2003, pp. 695–711.