

Proxemic Feature Recognition for Interactive Robots: Automating Metrics from the Social Sciences

Ross Mead, Amin Atrash, and Maja J. Matarić

Interaction Lab, Computer Science Department, University of Southern California
{rossmead, atrash, mataric}@usc.edu

Abstract. In this work, we discuss a set of metrics for analyzing human spatial behavior (proxemics) motivated by work in the social sciences. Specifically, we investigate individual, attentional, interpersonal, and physiological factors that contribute to social spacing. We demonstrate the feasibility of autonomous real-time annotation of these spatial features during multi-person social encounters. We utilize sensor suites that are non-invasive to participants, are readily deployable in a variety of environments (ranging from an instrumented workspace to a mobile robot platform), and do not interfere with the social interaction itself. Finally, we provide a discussion of the impact of these metrics and their utility in autonomous socially interactive systems.

Keywords: Proxemics, spatial interaction, spatial dynamics, social spacing, social robot, human-robot interaction, PrimeSensor, Microsoft Kinect.

1 Introduction and Background

Proxemics is the study of the dynamic process by which people position themselves in face-to-face social encounters [1]. This process is governed by sociocultural norms that, in effect, determine the overall sensory experience of each interacting participant [2]. People use proxemic signals, such as distance, stance, hip and shoulder orientation, head pose, and eye gaze, to communicate an interest in initiating, accepting, maintaining, terminating, or avoiding social interaction [3-5]. People can also manipulate space in an interaction, perhaps to direct attention to an external stimulus (usually accompanied by a hand gesture) or to guide a social partner to another location [6]. These cues are often subtle and noisy, and, subsequently, are subject to coarse analysis.

There exists a considerable body of work in the social sciences that seeks to analyze and explain certain proxemic phenomena. The anthropologist Edward T. Hall [1] coined the term “proxemics”, and proposed that *physiological* influences shaped by culture define zones of proxemic distances [2, 7]. Mehrabian [5], Argyle and Dean [8], and Burgoon *et al.* [9] analyzed spatial behaviors as a function of the *interpersonal* relationship between social partners. Schöne [10] was inspired by the spatial behaviors of biological organisms in response to stimuli, and investigated human spatial dynamics from physiological and *ethological* perspectives; similarly, Hayduk

and Mainprize [11] analyzed the personal space requirements of people who are blind. Kennedy *et al.* [12] studied the amygdala and how *emotional* (specifically, fight-or-flight) responses regulate space. Kendon [13] analyzed the *organizational* patterns of social encounters, categorizing them into *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.” Proxemic behavior is also impacted by factors of the *individual*—such as sex [14], age [15], ethnicity [16], and personality [17]—as well as *environmental* features—such as lighting [18], setting [19], location in setting and crowding [20], size [21], and permanence [7].

Following the emergence of embodied conversational agents [22], the study of proxemics was approached from a *computational* perspective. Rule-based social formation controllers have been applied to human-robot interaction (HRI) [23-25]. Interpersonal dynamic theories, such as equilibrium theory [8], have been evaluated in HRI [26, 27] and immersive virtual social environments [28]. Contemporary machine learning techniques have been applied to socially-appropriate person-aware navigation in dynamic crowded environments [29] and recognition of positive and negative attitudes of children with autism to an interactive robot [30].

A lack of high-resolution metrics limited previous efforts to coarse analyses in both space and time [31, 32]. Recent developments in markerless motion capture, such as the Microsoft Kinect¹ and the PrimeSensor², 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 human spatial interaction. In this work, we present a system that takes advantage of these advancements and draws on inspiration from existing metrics in the social sciences to automate the analysis process of proxemic behavior. This automation is necessary for the development of socially-situated artificial agents, both virtual and physically embodied.

2 Metrics for Proxemic Behavior Analysis

In this paper, we consider metrics that are commonly used by the social sciences to analyze proxemic behavior. Specifically, we are interested in the validated methods employed by Schegloff [33], McNeill [6, 34], Mehrabian [5], and Hall [35].

2.1 Schegloff’s Metrics

Schegloff [33] emphasized the importance of distinguishing between relative poses of the lower and upper parts of the body (Fig. 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”. He noted that, when a pose would deviate from its “home position” (i.e., 0°) with respect to an adjacent pose, the deviation would not last long and a compensatory orientation behavior

¹ <http://www.xbox.com/en-US/kinect>

² <http://www.primesense.com/>

would occur, 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 referred to this phenomenon as “body torque”.

- **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 to 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
- **Hip Torque:** angle between the hip and stance poses
- **Torso Torque:** angle between the torso and hip poses
- **Shoulder Torque:** angle between the shoulder and torso poses

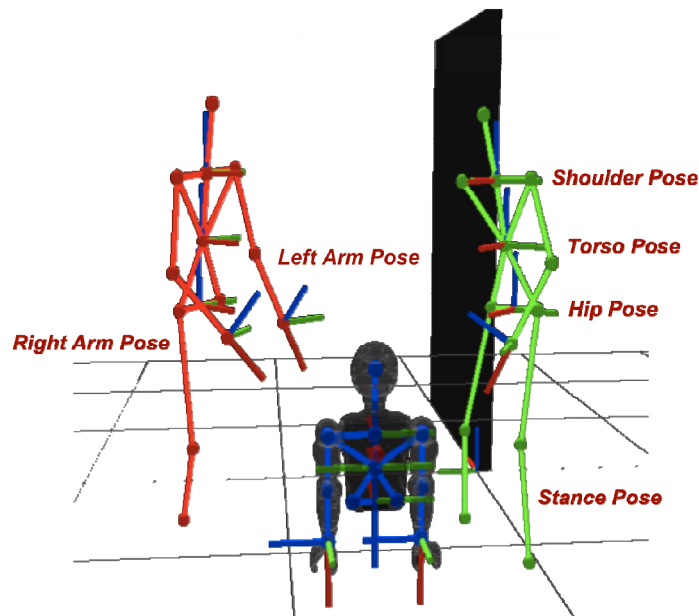


Fig. 1. Pose data for the two human users and an upper torso humanoid robot; the absence of some features, such as arms or legs, signified a position estimate with low confidence

2.2 McNeill's Metrics

McNeil [6, 34] posited that *attentional* (or *instrumental*) gestures could be used to manipulate space in social interactions by incorporating an object of mutual interest. As a first step to recognizing these features, we considered the relationship between the elbow and the hand of each individual in the interaction (Fig. 1).

- **Left Arm Pose:** position of left elbow, and orientation from left elbow to left hand
- **Right Arm Pose:** position of right elbow, and orientation from right elbow to right hand
- **Left Arm Deictic Referent:** target (social agent or stimulus) of a ray projected from the left elbow through the left hand
- **Right Arm Deictic Referent:** target (social agent or stimulus) of a ray projected from the right elbow through the right hand

2.3 Mehrabian's Metrics

Mehrabian [5] provides distance- and orientation-based interpersonal metrics for proxemic behavior analysis (Fig. 2). These spatial features are the most commonly used in the study of proxemics in both human-human and human-robot interactions.

- **Total Distance:** magnitude of a *Euclidean interpersonal distance vector* between the closest two points (one from each individual) in a social dyad
- **Straight-Ahead Distance:** magnitude of the *x*-component of the distance vector
- **Lateral Distance:** magnitude of the *y*-component of the distance vector
- **Relative Body Orientation:** magnitude of the angle between a social dyad

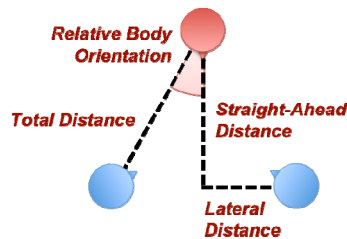


Fig. 2. In this triadic (three individual) interaction scenario, proxemic behavior is analyzed using simple spatial metrics between each social dyad (pair of individuals)

2.4 Hall's Metrics³

Hall's physiologically-inspired proxemic metrics are proposed as an alternative to strict spatial analysis, providing a sort of functional sensory explanation to the human use of space in social interaction [35] (Fig. 3). Hall seeks not only to answer questions of *where* and *when* a person will be, but, also, the question of *why*.

³For completeness, Hall's sex code (male or female) and posture code (prone, sitting, squatting, or standing) were also recorded by hand; automatic recognition of these proxemic factors was not attempted.

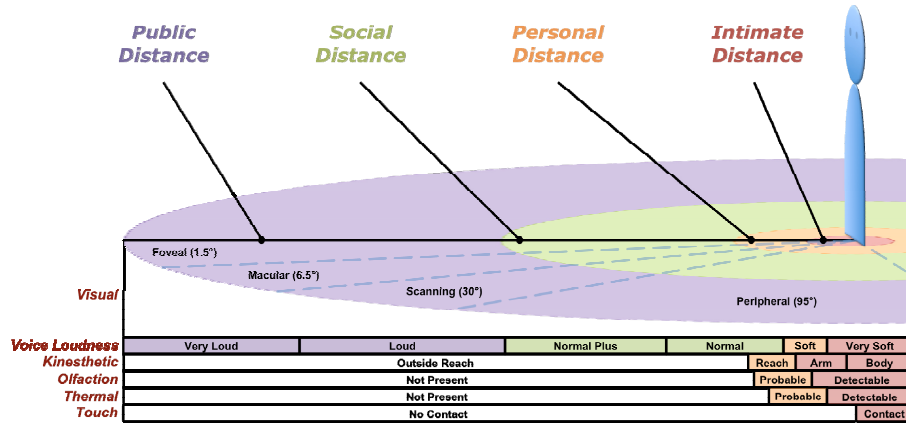


Fig. 3. Public, social, personal, and intimate distances, and the anticipated sensory sensations that an individual would experience while in each of these proximal zones

- **Distance Code**⁴: based on total distance; close-phase intimate distance (0"-6"), far-phase intimate distance (6"-18"), close-phase personal distance (18"-30"), far-phase personal distance (30"-48"), close-phase social distance (48"-84"), far-phase social distance (84"-144"), close-phase public distance (144"-300"), or far-phase public distance (more than 300")
- **Visual Code**: based on relative body orientation⁵; 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 arm poses; within body contact distance, just outside body contact distance, within easy touching distance with only forearm extended, just outside forearm distance ("elbow room"), within touching or grasping distance with the arms fully extended, just outside this distance, within reaching distance, or outside reaching distance
- **Olfaction Code**: based on total distance; differentiated body odor detectable (0"-6"), undifferentiated body odor detectable (6"-12"), breath detectable (12"-18"), olfaction probably present (18"-36"), or olfaction not present
- **Thermal Code**: based on total distance; conducted heat detected (0"-6"), radiant heat detected (6"-12"), heat probably detected (12"-21"), or heat not detected

⁴ These proxemic distances pertain to Western American culture—they are not cross-cultural.

⁵ In this implementation, shoulder pose was used to estimate the visual code, as the size of each person's face in the recorded image frames was too small to recognize eye gaze or head orientation.

- **Touch Code:** based on total distance; contact⁶ or no contact
- **Sociofugal-Sociopetal (SFP) Axis Code:** based on relative body orientation (in 22.5° intervals), with face-to-face (axis-0) representing maximum sociopetality and back-to-face (axis-8) representing maximum sociofugality [36-38]; axis-0 (0°-11.25°), axis-1 (11.25°-33.75°), axis-2 (33.75°-56.25°), axis-3 (56.25°-78.75°), axis-4 (78.75°-101.25°), axis-5 (101.25°-123.75°), axis-6 (123.75°-146.25°), axis-7 (146.25°-168.75°), or axis-8 (168.75°-180°)

3 Pilot Study

We conducted a pilot study to observe and analyze human spatial dynamics in natural interactions. The objective of this study was to demonstrate the feasibility of real-time annotation of proxemic behaviors in multi-person social encounters. In particular, we are interested in working with sensor suites that (1) are non-invasive to participants, (2) are readily deployable in a variety of environments (ranging from an instrumented workspace to a mobile robot), and (3) do not interfere with the interaction itself.

The study was set up and conducted in a twenty-by-twenty room in the Interaction Lab at the University of Southern California. A “presenter” and a participant engaged in a 5-6 minute open-ended (i.e., non-scripted) interaction loosely focused on a common object of interest—a static, non-interactive humanoid robot. (In this work, the second author played the role of the presenter, though was unaware of the conditions of the experiment.) The interactees were monitored by the PrimeSensor markerless motion capture system, an overhead color camera, and an omnidirectional microphone. A complete description of the experimental setup can be found in [39].

4 Qualitative Analysis and Discussion

A total of 18 participants were involved in the study. Joint positions recorded by the PrimeSensor were processed to recognize features based on the metrics of Schegloff [33], McNeill [6, 34], Mehrabian [5], and Hall [35], each with corresponding confidence values. We calculated these proxemic metrics and their respective changes for the participant, the presenter, and the robot in real-time. A preliminary validation of the feature recognition is available in [40]; a comprehensive analysis of system performance is currently underway.

For Schegloff’s metrics [33], pose information was accurately extracted from the PrimeSensor data [40]. Each joint position provided a confidence value, which was used to estimate the quality of the recognition of each individual pose. In addition, these poses simplified distance and orientation estimates for Mehrabian’s metrics [5].

To identify McNeill’s metrics [6], a ray was extended from the arm of the individual and used to detect collisions against any objects in the environment. This naïve

⁶ More formally, Hall’s touch code distinguishes between caressing and holding, feeling or caressing, extended or prolonged holding, holding, spot touching (hand peck), and accidental touching (brushing); however, automatic recognition of such forms of touching go beyond the scope of this work.

approach proved limited in detecting pointing behaviors in more complex social interactions—overall, the processed deictic referents were too noisy to be useful [40]. In future work, a model for recognizing pointing gestures will be trained to discriminate deictic behavior from other social and instrumental activities [34].

To our knowledge, this is the first time that Hall’s physiologically-inspired proxemic metrics [35] have been automatically extracted. By utilizing the distances between joint estimates, kinesthetic factors were accurately modeled for each individual [40]. Basic 3D collision detection techniques were used to identify contact between the participants (e.g., during a handshake). With the utilized sensor technologies, it would be feasible to consider a continuous measure of orientation, as opposed to the discretized SFP axis code. However, the sensors used in the current work were unable to provide eye gaze or head orientation; thus, metrics such as visual and olfactory receptor fields were estimated based on the recognized shoulder poses of the individuals. Future work will utilize additional sensors to estimate these metrics; for example, the Watson 3D head orientation tracker [41] could be used to estimate the field-of-view and olfactory fields, temperature and wind sensors to estimate thermal and olfactory factors, and more sensitive microphones to measure voice loudness. The addition of these sensors would provide the means to analyze spatial interaction dynamics in ways never before possible; improved recognition techniques via such sensors would improve their usability for autonomous systems. Finally, our anecdotal observations suggest that sensory (i.e., auditory, visual, thermal, olfactory) interference from the environment could extend the defined metrics and provide significant insight into the study of proxemic behavior [39]. In such cases, Hall’s metrics [35] might serve as better representations of the interaction than simple spatial metrics typically utilized in social encounters. Such an approach has yet to be investigated by any computational model, and will be a distinguishing factor in our continued work.

5 Conclusions and Future Work

In this work, we presented a set of metrics motivated by work in the social sciences for analyzing human spatial behavior. We then described a pilot study to demonstrate the feasibility of the real-time annotation of these spatial metrics in multi-person social encounters [39]. We utilized a markerless motion capture system (the Prime-Sensor), an overhead camera, and an omnidirectional microphone to recognize these metrics in real-time; this sensor suite is non-invasive to participants, is readily deployable in a variety of environments, and does not interfere with the interaction itself.

We are currently in the process of validating our automated system by comparing the performance against multiple human annotators; a preliminary validation is available in [40]. Our data set includes approximately 2 hours of data with 13 features per individual and 16 features per dyad. Although it is difficult to estimate the actual time for labeling by annotators, protocol within the USC Interaction Lab suggests allocating one hour per one minute of data for each feature; this further highlights the utility an automated annotation system. As future work focuses on collection of larger datasets, it becomes prohibitive to use human annotators to label the data. Additionally, as

autonomous agents use these metrics to conduct social interaction, it is important to understand the accuracy of the system for decision-making. An automated recognition system must provide sufficient accuracy to be useful.

We now present future directions and applications of this work, particularly with applications to autonomous systems.

5.1 Heuristic-Based vs. Learned Approaches

Currently, the models used by our automated system utilize heuristics based on metrics provided by the literature [5, 6, 33-35]. These parameters are often determined from empirical evidence as observed by the researchers. Consequently, this approach results in a discretization of the parameter space. However, advances in machine learning techniques allow for parameters to be determined automatically from data. This provides several advantages. First, it allows for continuous representation of the data. Current approaches discretize the parameter space into large regions. Automated systems can focus on learning continuous values of the data, resulting in a more fine-grained representation. Second, automatic learning allows for integration of large datasets; by taking advantage of automated labeling, models can be constructed from extensive datasets, which might be more accurate than hand-labeled models built from small datasets. Finally, learning of the models allows for impromptu and dynamic scenarios. These parameters may vary from situation-to-situation and person-to-person. By adjusting the parameters, an agent could adapt to each individual.

5.2 Environmental Interference

The study conducted in this work was performed in a controlled idealized social scenario; however, real-world scenarios will include multiple sources of environmental interference, including noise from a crowded room and/or low-light conditions [18]. An autonomous agent might encounter such dynamic and unstructured factors in social encounters; thus, future work will examine the use of the system in these real-world scenarios, examining its ability to recognize the metrics described, as well as exploring additional techniques to handle potential errors. In particular, we will focus on the impact of such interference on Hall's metrics [35].

5.3 Human-Robot Interaction

As autonomous intelligent robots become more ubiquitous in society, it will become necessary for them to interact with humans using natural modalities. Recognition of spatial behaviors, such as those described in this work, is a fundamental component for these agents to interact in an appropriate manner. By combining this study of proxemics with techniques from robotics and machine learning, we can create robust, socially-aware robots. Sensor data fed into a predictive model would allow the agent to recognize high-level interaction cues. A generative model incorporated into an action selection system would allow the agent to produce the similar interaction cues as it was trained to recognize, in an effort to produce a predictable social response from the human.

Acknowledgements. This work is supported in part by an NSF Graduate Research Fellowship, as well as ONR MURI N00014-09-1-1031 and NSF CNS-0709296 grants, and the Dan Marino Foundation. We would like to thank Louis-Philippe Morency for his insight in the experimental design process, and Mark Bolas and Evan Suma for their assistance in using the PrimeSensor. We would also like to thank Edward Kaszubski for his help in constructing the physical setup, and Dave Feil-Seifer for his assistance in using the overhead camera and Overhead Interaction Toolkit.

References

1. Hall, E.T.: *The Silent Language*. Doubleday Co., New York (1959)
2. Hall, E.T.: *Handbook for Proxemic Research*. American Anthropology Assn., Washington (1974)
3. Deutsch, R.D.: *Spatial Structurings in Everyday Face-to-face Behavior: A Neurocybernetic Model*. In: *The Association for the Study of Man-Environment Relations*, Orangeburg (1977)
4. Mead, R., Atrash, A., Matarić, M.J.: *Recognition of Spatial Dynamics for Predicting Social Interaction*. In: *HRI 2011*, Lausanne (2011)
5. Mehrabian, A.: *Nonverbal Communication*. Aldine Transaction, Piscataway (1972)
6. McNeill, D.: *Gesture and Thought*. Chicago University Press, Chicago (2005)
7. Hall, E.T.: *The Hidden Dimension*. Doubleday Co., New York (1966)
8. Argyle, M., Dean, J.: *Eye-Contact, Distance, and Affiliation*. *Sociometry* 28, 289–304 (1965)
9. Burgoon, J.K., Stern, L.A., Dillman, L.: *Interpersonal Adaptation: Dyadic Interaction Patterns*. Cambridge University Press, New York (1995)
10. Schöne, H.: *Spatial Orientation: The Spatial Control of Behavior in Animals and Man*. Princeton University Press, Princeton (1984)
11. Hayduk, L.A., Mainprize, S.: *Personal Space of the Blind*. *Social Psychology Quarterly* 43(2), 216–223 (1980)
12. Kennedy, D.P., Gläscher, J., Tyszka, J.M., Adolphs, R.: *Personal Space Regulation by the Human Amygdala*. *Nature Neuroscience* 12, 1226–1227 (2009)
13. Kendon, A.: *Conducting Interaction - Patterns of Behavior in Focused Encounters*. Cambridge University Press, New York (1990)
14. Price, G.H., Dabbs Jr., J.M.: *Sex, Setting, and Personal Space: Changes as Children Grow Older*. *Personal Social Psychology Bulletin* 1, 362–363 (1974)
15. Aiello, J.R., Aiello, T.D.: *The Development of Personal Space: Proxemic Behavior of Children 6 through 16*. *Human Ecology* 2(3), 177–189 (1974)
16. Jones, S.E., Aiello, J.R.: *Proxemic Behavior of Black and White First-, Third-, and Fifth-Grade Children*. *Journal of Personality and Social Psychology* 25(1), 21–27 (1973)
17. Aiello, J.R.: *Human Spatial Behavior*. In: *Handbook of Environmental Psychology*, ch. 12. John Wiley & Sons, New York (1987)
18. Adams, L., Zuckerman, D.: *The Effect of Lighting Conditions on Personal Space Requirements*. *The Journal of General Psychology* 118(4), 335–340 (1991)
19. Geden, E.A., Begeman, A.V.: *Personal Space Preferences of Hospitalized Adults*. *Research in Nursing and Health* 4, 237–241 (1981)
20. Evans, G.W., Wener, R.E.: *Crowding and Personal Space Invasion on the Train: Please Don't Make Me Sit in the Middle*. *Journal of Environmental Psychology* 27, 90–94 (2007)

21. Aiello, J.R., Thompson, D.E., Brodzinsky, D.M.: 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)
22. Cassell, J., Sullivan, J., Prevost, S., Churchill, E.F.: *Embodied Conversational Agents*. MIT Press, Cambridge (2000)
23. Kuzuoka, H., Suzuki, Y., Yamashita, J., Yamazaki, K.: Reconfiguring Spatial Formation Arrangement by Robot Body Orientation. In: *HRI 2010, Osaka*, pp. 285–292 (2010)
24. Walters, M.L., Dautenhahn, K., Boekhorst, R.T., Koay, K.L., Syrdal, D.S., Nehaniv, C.L.: An Empirical Framework for Human-Robot Proxemics. In: *New Frontiers in Human-Robot Interaction, Edinburgh*, pp. 144–149 (2009)
25. Huettenrauch, H., Eklundh, K.S., Green, A., Topp, E.A.: Investigating Spatial Relationships in Human-Robot Interaction. In: *IROS 2006, Beijing* (2006)
26. Mumm, J., Mutlu, B.: Human-Robot Proxemics: Physical and Psychological Distancing in Human-Robot Interaction. In: *HRI 2011, Lausanne*, pp. 331–338 (2011)
27. Takayama, L., Pantofaru, C.: Influences on Proxemic Behaviors in Human-Robot Interaction. In: *IROS 2009, St. Louis* (2009)
28. Bailenson, J.N., Blascovich, J., Beall, A.C., Loomis, J.M.: Equilibrium Theory Revisited: Mutual Gaze and Personal Space in Virtual Environments. *Presence* 10(6), 583–598 (2001)
29. Trautman, P., Krause, A.: Unfreezing the Robot: Navigation in Dense, Interacting Crowds. In: *IROS 2010, Taipei* (2010)
30. Feil-Seifer, D.J., Matarić, M.J.: Automated Detection and Classification of Positive vs. Negative Robot Interactions with Children with Autism using Distance-based Features. In: *HRI 2011, Lausanne*, pp. 323–330 (2011)
31. Oosterhout, T.V., Visser, A.: A Visual Method for Robot Proxemics Measurements. In: *HRI 2008 Workshop on Metrics for Human-Robot Interaction, Technical Report 471, Amsterdam* (2008)
32. Jones, S.E., Aiello, J.R.: A Test of the Validity of Projective and Quasi-projective measures of Interpersonal Distance. *The Western Journal of Speech Communication* 43, 143–152 (1979)
33. Schegloff, E.A.: Body Torque. *Social Research* 65(3), 535–596 (1998)
34. McNeill, D.: *Gesture, Gaze, and Ground*. LNCS, pp. 1–14. Springer, Heidelberg (2005)
35. Hall, E.T.: A System for Notation of Proxemic Behavior. *American Anthropologist* 65, 1003–1026 (1963)
36. Sommer, R.: Sociofugal Space. *The American Journal of Sociology* 72(6), 654–660 (1967)
37. Lawson, B.: *Sociofugal and Sociopetal Space, The Language of Space*. Architectural Press, Oxford (2001)
38. Low, S.M., Lawrence-Zúñiga, D.: *The Anthropology of Space and Place: Locating Culture*. Blackwell Publishing, Oxford (2003)
39. Mead, R., Matarić, M.J.: An Experimental Design for Studying Proxemic Behavior in Human-Robot Interaction. Technical Report CRES-11-001, USC Interaction Lab, Los Angeles (2011)
40. Mead, R., Atrash, A., Matarić, M.J.: A Preliminary Validation of Automated Feature Recognition for Proxemic Behavior Analysis. Technical Report CRES-11-002, USC Interaction Lab, Los Angeles (2011)
41. Morency, L.P., Sidner, C., Lee, C., Darrell, T.: Contextual Recognition of Head Gestures. In: *International Conference on Multimodal Interactions, Toronto*, pp. 18–24 (2005)