

Automated Proxemic Feature Extraction and Behavior Recognition: Applications in Human-Robot Interaction

Ross Mead · Amin Atrash · Maja J. Matarić

Accepted: 2 May 2013
© Springer Science+Business Media Dordrecht 2013

Abstract In this work, we discuss a set of feature representations for analyzing human spatial behavior (proxemics) motivated by metrics used in the social sciences. Specifically, we consider individual, physical, and psychophysical factors that contribute to social spacing. We demonstrate the feasibility of autonomous real-time annotation of these proxemic features during a social interaction between two people and a humanoid robot in the presence of a visual obstruction (a physical barrier). We then use two different feature representations—physical and psychophysical—to train Hidden Markov Models (HMMs) to recognize spatiotemporal behaviors that signify transitions into (initiation) and out of (termination) a social interaction. We demonstrate that the HMMs trained on psychophysical features, which encode the sensory experience of each interacting agent, outperform those trained on physical features, which only encode spatial relationships. These results suggest a more powerful representation of proxemic behavior with particular implications in autonomous socially interactive and socially assistive robotics.

Keywords Proxemics · Spatial interaction · Spatial dynamics · Sociable spacing · Social robot · Human-robot interaction · PrimeSensor · Microsoft Kinect

R. Mead (✉) · A. Atrash · M.J. Matarić
Interaction Lab, University of Southern California, Los Angeles,
USA
e-mail: rossmead@usc.edu

A. Atrash
e-mail: atrash@usc.edu

M.J. Matarić
e-mail: mataric@usc.edu

1 Introduction

Proxemics is the study of the dynamic process by which people position themselves in face-to-face social encounters [14]. This process is governed by sociocultural norms that, in effect, determine the overall sensory experience of each interacting participant [17]. 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 [10, 35, 36]. These cues are often subtle and noisy.

A lack of high-resolution metrics limited previous efforts to model proxemic behavior to coarse analyses in both space and time [24, 39]. Recent developments in markerless motion capture—such as the Microsoft Kinect,¹ the Asus Xtion,² and the PrimeSensor³—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 (features) in the social sciences to automate the analysis of proxemics. We then utilize these extracted features to recognize spatiotemporal behaviors that signify transitions into (initiation) and out of (termination) a social interaction. This automation is necessary for the development of socially situated artificial agents, both virtually and physically embodied.

¹<http://www.xbox.com/kinect>.

²http://www.asus.com/Multimedia/Motion_Sensor.

³<http://www.primesense.com>.

2 Background

While the study of proxemics in human-agent interaction is relatively new, there exists a rich body of work in the social sciences that seeks to analyze and explain proxemic phenomena. Some proxemic models from the social sciences have been validated in human-machine interactions. In this paper, we focus on models from the literature that can be applied to recognition and control of proxemics in human-robot interaction (HRI).

2.1 Proxemics in the Social Sciences

The anthropologist Edward T. Hall [14] coined the term “proxemics”, and proposed that *psychophysical* influences shaped by culture define zones of proxemic distances [15–17]. Mehrabian [36], Argyle and Dean [5], and Burgoon et al. [7] analyzed *psychological* indicators of the interpersonal relationship between social dyads. Schöne [45] 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 [18] analyzed the personal space requirements of people who are blind. Kennedy et al. [27] studied the amygdala and how *emotional* (specifically, fight-or-flight) responses regulate space. Kendon [26] 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.” Proxemics is also impacted by factors of the *individual*—such as involvement [44], sex [41], age [3], ethnicity [23], and personality [2]—as well as *environmental* features—such as lighting [1], setting [13], location in setting and crowding [11], size [4], and permanence [16].

2.2 Proxemics in Human-Agent Interaction

The emergence of embodied conversational agents (ECAs) [8] necessitated *computational* models of social proxemics. Proxemics for ECAs has been parameterized on culture [22] and so-called “social forces” [21]. The equilibrium theory of spacing [5] was found to be consistent with ECAs in [6, 31]. Pelachaud and Poggi [40] provides a summary of aspects of emotion, personality, and multimodal communication (including proxemics) that contribute to believable embodied agents.

As the trend in robotic systems places them in social proximity of a human user, it becomes necessary for them to interact with humans using natural modalities. Understanding social spatial behaviors is a fundamental concept for these agents to interact in an appropriate manner. Rule-based proxemic controllers have been applied to HRI [20,

29, 52]. Interpersonal dynamic models, such as [5], have been investigated in HRI [38, 48]. Jung et al. [25] proposed guidelines for robot navigation in speech-based and speechless social situations (related to Hall’s [15] *voice loudness code*), targeted at maximizing user acceptance and comfort. A spatially situated methodology for evaluating “interaction quality” (social presence) in mobile remote telepresence interactions has been proposed and validated based on Kendon’s [26] theory of proxemic F-formations [28]. Contemporary probabilistic modeling techniques have been applied to socially appropriate person-aware robot navigation in dynamic crowded environments [50, 51], to calculate a robot approach trajectory to initiate interaction with a walking person [43], to recognize the averse and non-averse reactions of children with autism spectrum disorder to a socially assistive robot [12], and to position the robot for user comfort [49]. A lack of high-resolution quantitative measures has limited these efforts to coarse analyses [23, 39].

In this work, we present a set of feature representations for analyzing proxemic behavior motivated by metrics commonly used in the social sciences. Specifically, in Sect. 3, we consider individual, physical, and psychophysical factors that contribute to social spacing. In Sects. 4 and 5, we demonstrate the feasibility of autonomous real-time annotation of these proxemic features during a social interaction between two people and a humanoid robot. In Sect. 6, we compare the performance of two different feature representations to recognize high-level spatiotemporal social behaviors. In Sect. 7, we discuss the implications of these representations of proxemic behavior, and propose an extension for more continuous and situated proxemics in HRI [34].

3 Feature Representation and Extraction

In this work, proxemic “features” are based on the most commonly used metrics in the social sciences literature. We first extract Schegloff’s *individual* features [44]. We then use the features of two individuals to calculate the features for the dyad (pair). We are interested in two popular and validated annotation schemas of proxemics: (1) Mehrabian’s *physical* features [36], and (2) Hall’s *psychophysical* features [15].

3.1 Individual Features

Schegloff [44] 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 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 [44] referred to this phenomenon as *body torque*, which has been investigated in HRI [29].

We used the following list of individual features:

- **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; extracted and tracked using the head pose estimation technique described in [37]
- **Hip Torque:** angle between hip and stance poses
- **Torso Torque:** angle between torso and hip poses
- **Shoulder Torque:** angle between shoulder and torso poses
- **Head Torque:** angle between head and shoulder poses.

3.2 Physical Features

Mehrabian [36] provides distance- and orientation-based metrics between a dyad (two individuals) for proxemic behavior analysis (Fig. 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 social dyad between agents *A* and *B*:

- **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*

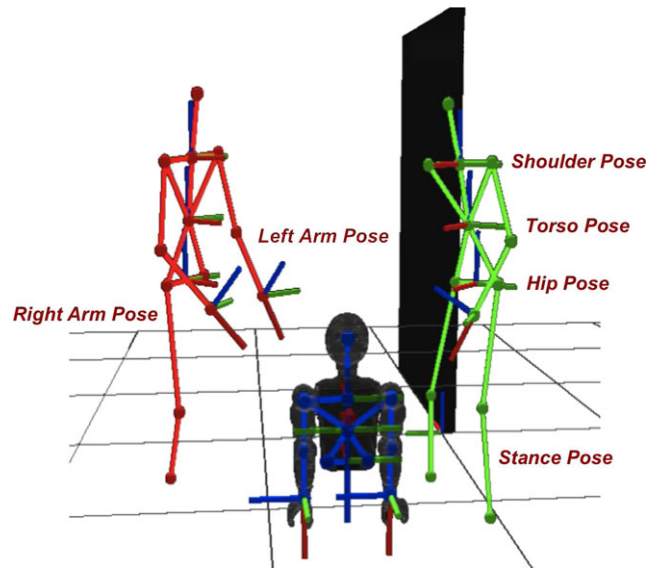


Fig. 1 Pose data for two human users and an upper torso humanoid robot; the absence of some features—such as head, arms, or legs—signified a pose estimate with low confidence

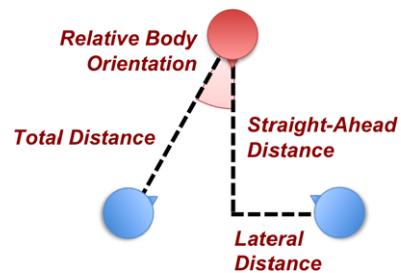


Fig. 2 In this triadic (three agent) interaction scenario, proxemic behavior is analyzed using simple physical metrics between each social dyad (pair of individuals)

3.3 Psychophysical Features

Hall's [15] *psychophysical* proxemic metrics are proposed as an alternative to strict physical analysis, providing a sort of functional sensory explanation to the human use of space in social interaction (Fig. 3). Hall [15] seeks not only to answer questions of *where* a person will be, but, also, the question of *why* they are there, investigating the underlying processes and systems that govern proxemic behavior.

For example, upon first meeting, two Western American strangers often shake hands, and, in doing so, subconsciously gauge 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 [19]. 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

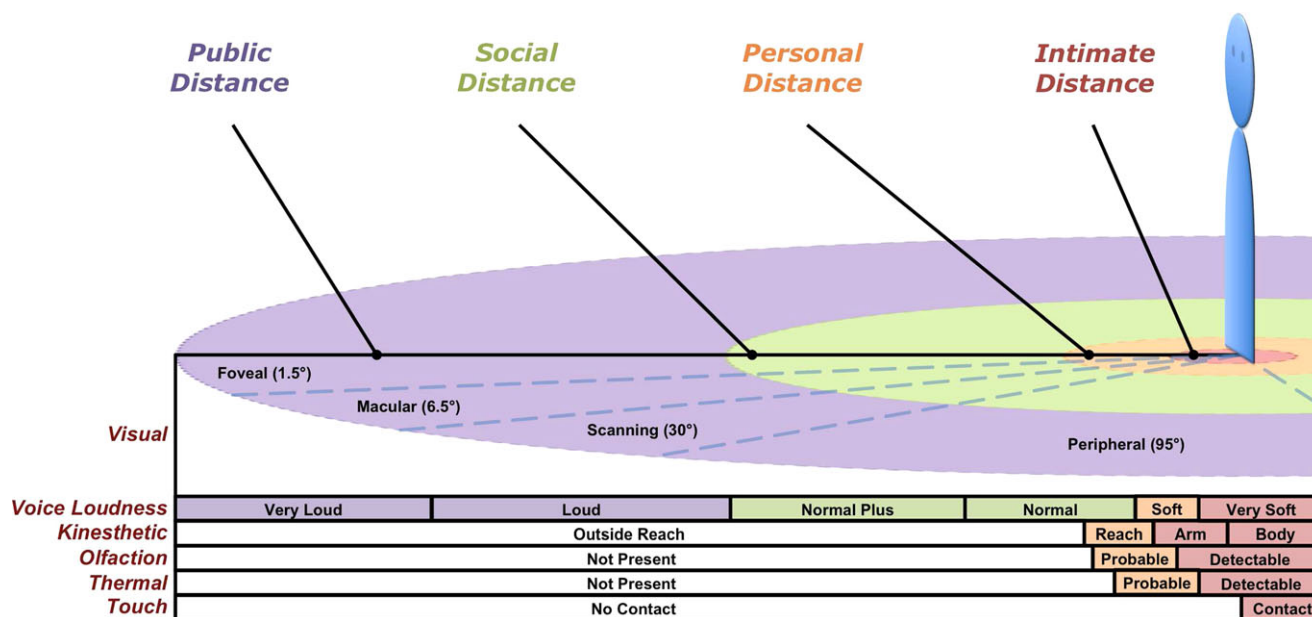


Fig. 3 Public, social, personal, and intimate distances, and the anticipated sensory sensations that an individual would experience while in each of these proximal zones

the fact that brief physical embrace (e.g., hugging) is common at this range [15]. However, olfactory and thermal sensations of one another 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, 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” [15].

Hall’s [15] coding schema is typically annotated by social scientists based purely on distance and orientation data observed from video [16]. The automation of this tedious process is a major contribution of this work; to our knowledge, this is the first time that these proxemic features have been automatically extracted.

The psychophysical “feature codes” and their corresponding “feature intervals” for each individual in a social dyad between agents *A* and *B* are as follows:

- **Distance Code:**⁴ 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 [30,

32, 47]; 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°)

- **Visual Code:** based on head pose;⁵ 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

⁵In this implementation, head pose was used to estimate the visual code; however, as the size of each person’s face in the recorded image frames was rather small, the results from the head tracker were quite noisy [37]. If the head pose estimation confidence was below some threshold, the system would instead rely on the shoulder pose for eye gaze estimates.

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

- **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
- **Touch Code:** based on total distance;⁶ contact⁷ or no contact

4 System Implementation and Discussion

The feature extraction system can be implemented using any human motion capture technique. We utilized the PrimeSensor structured light range sensor and the OpenNI⁸ person tracker for markerless motion capture. We chose this setup because it (1) is non-invasive to participants, (2) is readily deployable in a variety of environments (ranging from an instrumented workspace to a mobile robot), and (3) 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). We developed an error model of the sensor, which was then used to generate error models for individual [44], physical [36], and psychophysical [15] features, discussed below.

4.1 Motion Capture System

We conducted an evaluation of the precision of the PrimeSensor distance estimates. The PrimeSensor was mounted atop a tripod at a height of 1.5 meters and pointed straight at a wall. We placed the sensor rig at 0.2-meter intervals between 0.5 meters and 2.5 meters and, at each location, took distance readings (a collection of 3-dimensional points, referred to as a “point cloud”) at each location. We used a planar model segmentation technique⁹ to eliminate points in the point cloud that did not fit onto the wall plane. We calculated the average depth reading of each point in the segmented plane, and modeled the sensor error E as a function of distance d (in meters): $E(d) = k \times d^2$, with $k = 0.0055$. Our procedure and results are consistent with that reported in the

⁶In this implementation, we utilized the 3-dimensional point cloud provided by our motion capture system for improved accuracy (see Sect. 4.1); however, we assume nothing about the implementations of others, so total distance can be used for approximations in the general case.

⁷More 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 extraction of such forms of touching go beyond the scope of this work.

⁸<http://openni.org>.

⁹http://www.pointclouds.org/documentation/tutorials/planar_segmentation.php.

ROS Kinect accuracy test,¹⁰ as well as those of other structured light and stereo camera range estimates, 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 .

4.2 Individual Features

Shotton et al. [46] provides a comprehensive evaluation of individual joint pose 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, 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 may differ, the performance is comparable for our purposes.

4.3 Physical Features

In our estimation of 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 [46], this value is learned, with an average offset of 0.039 meters. To extract accurate physical proxemic features of the social dyad, we subtract twice this value (once for each individual) 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 OpenNI software is beyond the scope and resources of this work; instead, we refer to the comparable figures reported in [46].

4.4 Psychophysical Features

Each feature annotation in Hall’s [15] psychophysical representation was developed based upon values from literature on the human sensory system [16]. 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 interval¹¹—such evaluations come from literature cited by Hall [15, 16] when the representation was initially proposed. Rather, this work 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.

Intervals for each feature code were evaluated at 1, 2, 3, and 4 meters from the sensor; the sensor was orthogonal to

¹⁰http://www.ros.org/wiki/openni_kinect/kinect_accuracy.

¹¹For example, we do not measure the radiant heat or odor transmitted by one individual and the intensity at the corresponding sensory organ of the receiving individual.

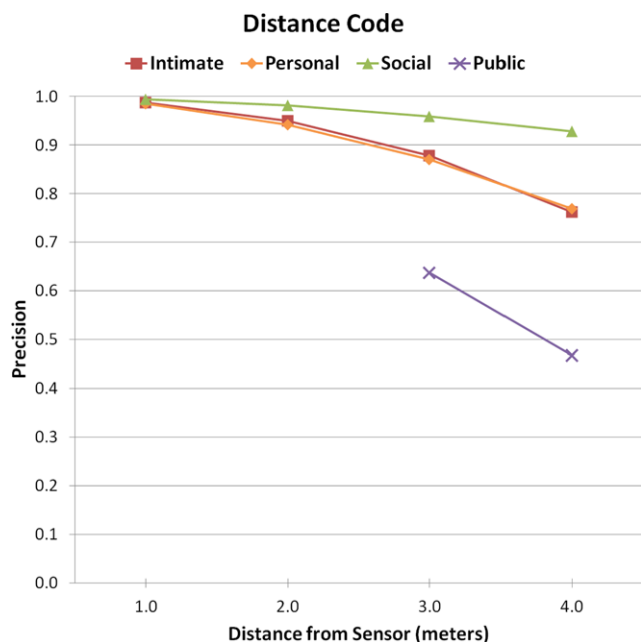


Fig. 4 Error model of psychophysical distance code

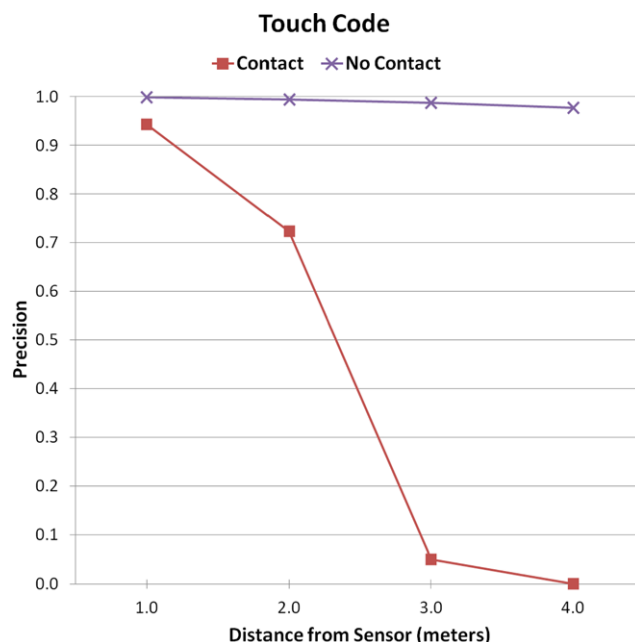


Fig. 6 Error model of psychophysical touch code

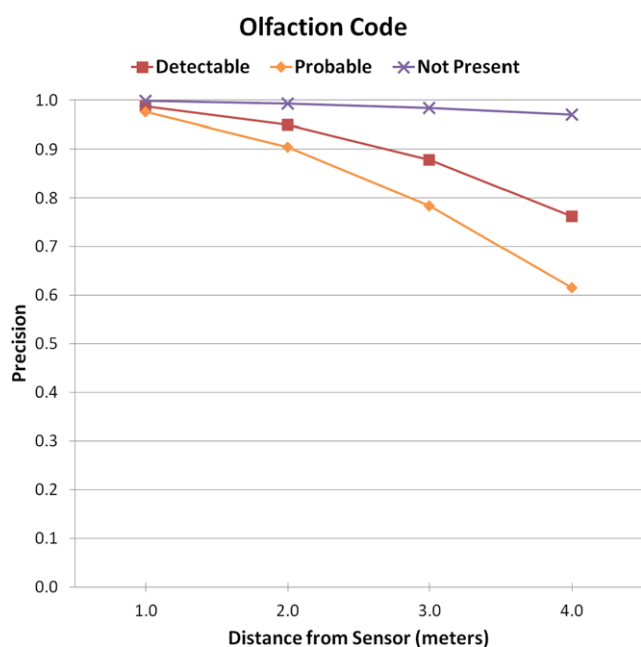


Fig. 5 Error model of psychophysical olfaction code

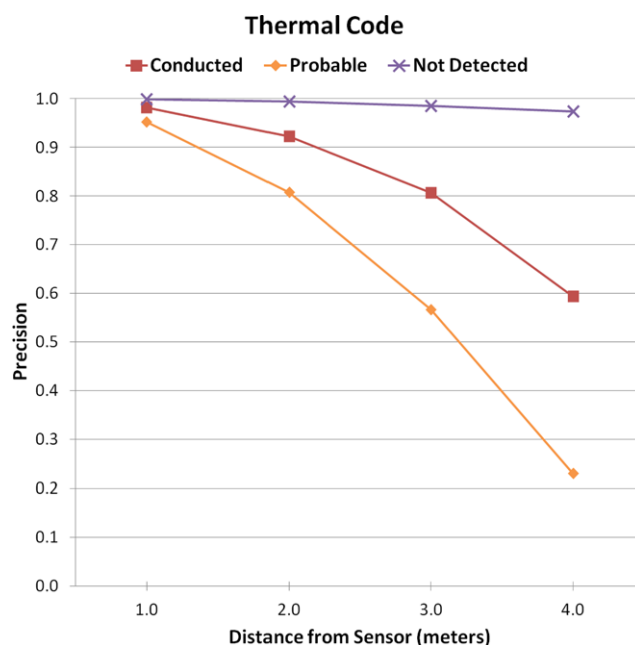


Fig. 7 Error model of psychophysical thermal code

a line passing through the hip poses of two standing individuals. The results of these estimates are illustrated in Figs. 4–11. At some ranges from the sensor, a feature interval would place the individuals out of the sensor field-of-view; these values are omitted.¹²

¹²This occurs for the SFP axis estimates, as well as at the public distance interval, the loud and very loud voice loudness intervals, and the outside reach kinesthetic interval.

The *SFP axis code* contains feature intervals of uniform size (40°). Rather than evaluate each interval independently, we evaluated the average precision of the intervals at different distances. We considered the error in estimated orientation of one individual *w.r.t.* another at 1, 2, 3, and 4 meters from each other. At shorter ranges, the error in estimated position of one increases the uncertainty of the orientation estimate (Fig. 9).

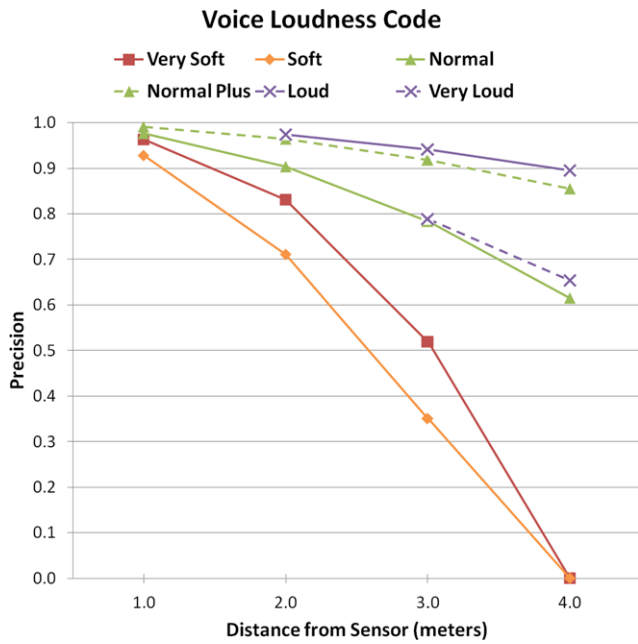


Fig. 8 Error model of psychophysical voice loudness code

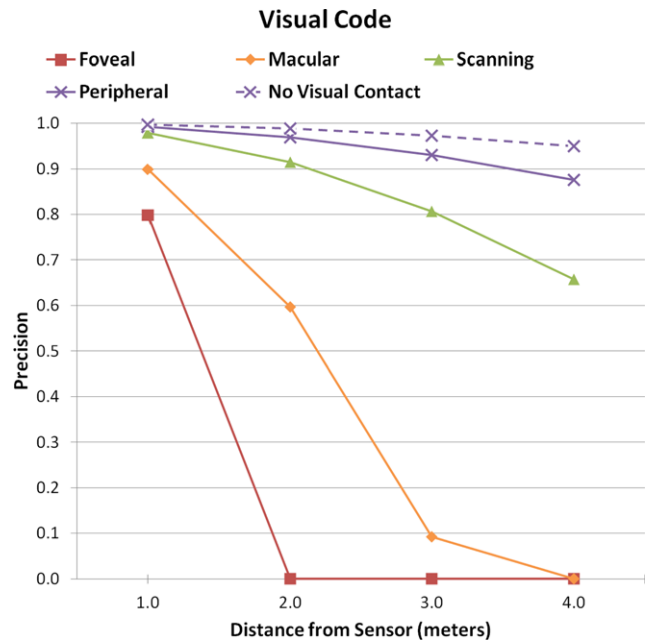


Fig. 10 Error model of psychophysical visual code

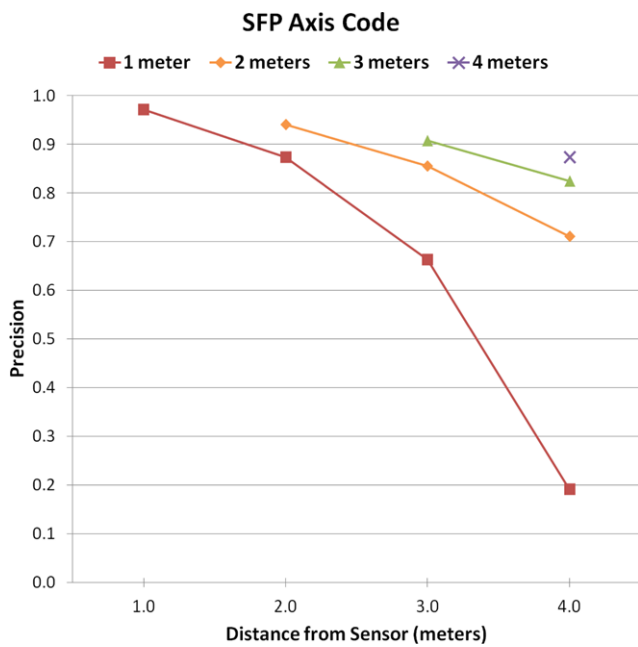


Fig. 9 Error model of psychophysical sociofugal-sociopetal (SFP) axis code

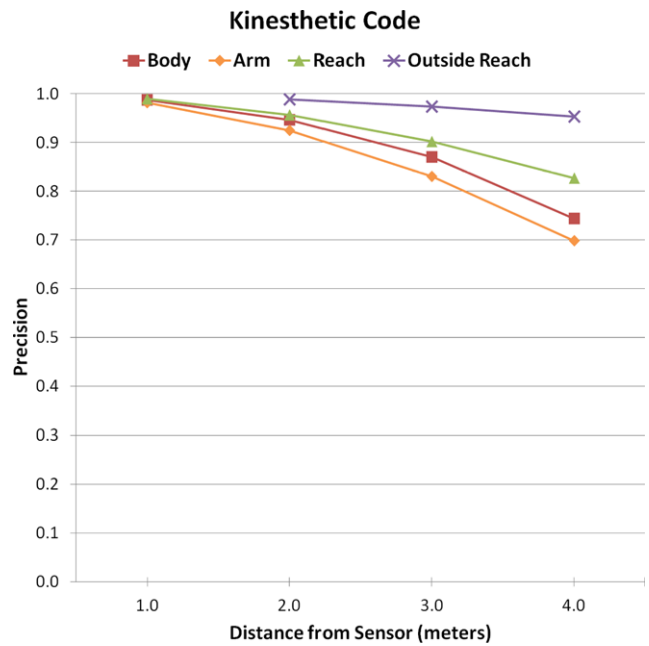


Fig. 11 Error model of psychophysical kinesthetic code

Eye gaze was unavailable for true *visual code* extraction. Instead, the *visual code* was estimated based on the head pose [37]; when this approach failed, shoulder pose would be used to estimate eye gaze. Both of these estimators resulted in coarse estimates of narrow feature intervals (foveal and macular) (Fig. 10). We are collaborating with the researchers of [37] to improve the performance of the head

pose estimation system and, thus, the estimation of the *visual code*, at longer ranges.

For evaluation of the *kinesthetic code*, we used interval distance estimates based on average human limb lengths (Fig. 11) [16]. 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 is expected to be more misclassifications at the feature interval boundaries.

5 Interaction Study

We conducted an interaction study to observe and analyze human proxemic behavior in natural social encounters [33]. The study was approved by the Institutional Review Board #UP-09-00204 at the University of Southern California. Each participant was presented with the Experimental Subject's Bill of Rights, and was informed of the general nature of the study and the types of data that were being captured (video and audio).

The study objectives, setup, and interaction scenario (context) are summarized below. A full description of the interaction study can be found in [33].

5.1 Objectives

The objective of this study was to demonstrate the utility of real-time annotation of proxemic features to recognize higher-order spatiotemporal behaviors in multi-person social encounters. To do this, we sought to capture proxemic behaviors signifying transitions into (initiation) and out of (termination) social interactions [10, 35]. *Initiation behavior* attempts to engage or recognize a potential social partner in discourse (also referred to as a “sociopetal” behavior [30, 32]). *Termination behavior* proposes the end of an interaction in a socially appropriate manner (also referred to as a “sociofugal” behavior [30, 32]). These behaviors are directed at a social stimulus (i.e., an object or another agent), and occur sequentially or in parallel *w.r.t.* each stimulus.

5.2 Setup

The study was set up and conducted in a 20'-by-20' room in the Interaction Lab at the University of Southern California (Fig. 12). A “presenter” and a participant engaged in an interaction loosely focused on a common object of interest—a static, non-interactive humanoid robot. The interactees were monitored by the PrimeSensor markerless motion capture system, an overhead color camera, and an omnidirectional microphone.

Prior to the participant entering the room, the presenter stood on floor marks X and Y for user calibration. The participant later entered the room from floor mark A, and awaited sensor calibration at floor marks B and C; note that, from all participant locations, the physical divider obstructed the participant's view of the presenter (i.e., the participant could not see and was not aware that the presenter was in the room).

A complete description of the experimental setup and data collection systems can be found in [33].

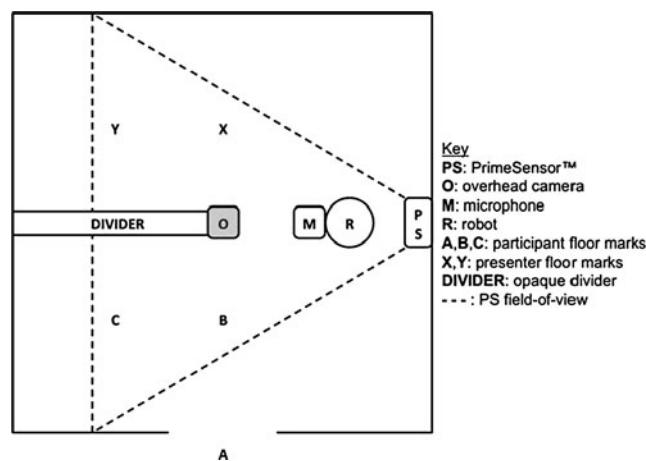


Fig. 12 The experimental setup

5.3 Scenario

As soon as the participant moved away from floor mark C and approached the robot (an initiation behavior directed at the robot), the scenario was considered to have officially begun. Once the participant verbally engaged the robot (unaware that the robot would not respond), the presenter was signaled (via laser pointer out of the field-of-view of the participant) to approach the participant from behind the divider, and attempt to enter the existing interaction between the participant and the robot (an initiation behavior directed at both the participant and the robot, often eliciting an initiation behavior from the participant directed at the presenter). Once engaged in this interaction, the dialogue between the presenter and the participant was open-ended (i.e., unscripted) and lasted 5–6 minutes. Once the interaction was over, the participant exited the room (a termination behavior directed at both the presenter and the robot); the presenter had been previously instructed to return to floor mark Y at the end of the interaction (a termination behavior directed at the robot). Once the presenter reached this destination, the scenario was considered to be complete.

5.4 Dataset

A total of 18 participants were involved in the study. Joint positions recorded by the PrimeSensor were processed to extract individual [44], physical [36], and psychophysical [15] features discussed in Sect. 3.

The data collected from these interactions were annotated with the behavioral events *initiation* and *termination* based on each interaction dyad (i.e., behavior of one social agent A directed at another social agent B). The dataset provided 71 examples of initiation and 69 examples of termination. Two sets of features were considered for comparison: (a) Mehrabian's [36] physical features, capturing distance and orientation; and (b) Hall's [15] psychophysical features,

capturing the sensory experience of each agent. Physical features included total, straight-ahead, and lateral distances, as well as orientation of agent *B* with respect to agent *A* [36]; similarly, psychophysical features included SFP axis, visual, voice loudness, kinesthetic, olfaction, thermal, and touch codes [15] (Fig. 14). All of these features were automatically extracted using the system described above.

6 Behavior Modeling and Recognition

To examine the utility of these proxemic feature representations, data collected from the pilot study were used to train an automated recognition system for detecting social events, specifically, *initiation* and *termination* behaviors (see Sect. 5.1).

6.1 Hidden Markov Models

Hidden Markov Models (HMMs) are stochastic processes for modeling of nondeterministic time sequence data [42]. HMMs are defined by a set of *states*, $s \in S$, and a set of *observations*, $o \in O$. A *transition function*, $P(s'|s)$, defines the likelihood of transitioning from state s at time t to state s' at time $t + 1$. When entering a state, s , the agent then receives an observation based on the distribution $P(o|s)$. These models rely on the Markov assumption, which states that the conditional probability of future states and observations depend only on the current state of the agent. HMMs are commonly used in recognition of time-sequence data, such as speech, gesture, and behavior [42].

For multidimensional data, the observation is often factored into independent features, f_i , and treated as a vector, where $o = f_1, f_2, \dots, f_n$. The resulting likelihood is the joint distribution of all of the likelihood of the features, such that $P(o|s) = \prod_{i=1..n} P(f_i|s)$. For continuous features, $P(f_i|s)$ is maintained as a Gaussian distribution defined by a mean and variance.

When used for discrimination between multiple classes, a separate model, M_j , is trained for each class, j (e.g., initiation and termination behaviors). For a given observation sequence, o_1, o_2, \dots, o_m , the likelihood of that sequence given each model, $P(o_1, o_2, \dots, o_m|M_j)$, is calculated. The sequence is then labeled (classified) as the most likely model. Baum-Welch is an expectation-maximization algorithm used to train the parameters of the models given a data set [9]. For a data set, Baum-Welch determines the parameters that maximize the likelihood of the data occurring.

In this work, we utilized five-state left-right HMMs with two skip-states (Fig. 13); this is a common topology used in recognition applications [42]. Observation vectors for each representation (physical and psychophysical) consisted of 7 features and 11 features, respectively (Fig. 14). For each representation, two HMMs were trained: one for initiation and

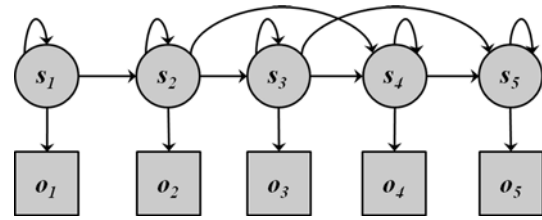


Fig. 13 A five-state left-right HMM with two skip-states [42] used to model each behavior, M_j , for each representation

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

Fig. 14 The observation vectors, o_i , for 7-dimensional physical features (left) and 11-dimensional psychophysical features (right)

one for termination. When a new behavior instance was observed, the models returned the likelihood of that instance being initiation or termination. The observation and transition parameters were processed and converged after six iterations of the Baum-Welch algorithm [9]. Leave-one-out cross-validation was utilized to validate the performance of the models.

6.2 Results and Analysis

Table 1 presents the results when training the models using the physical features. In this case, while the system is able to discriminate between the two models, there is often misclassification, resulting in an overall accuracy of 56 % (Fig. 15). 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).

Table 2 presents results using the psychophysical features, showing considerable improvement, with an overall accuracy of 72 % (Fig. 15). Psychophysical features attempt to account for an agent’s sensory experience 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, we anticipate that the relative performance between them would remain unchanged.

The intuition behind the improvement in performance of the psychophysical representation over the physical repre-

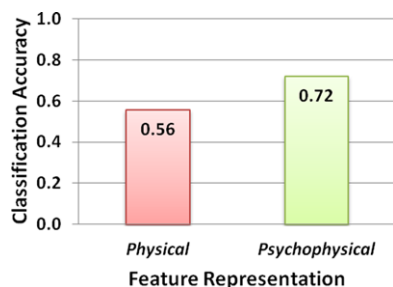


Fig. 15 Comparison of HMM classification accuracy of initiation and termination behaviors trained over physical and psychophysical feature sets

Table 1 Confusion matrix for recognizing initiation and termination behaviors using physical features

Example / Classification	Initiation	Termination
Initiation	41	30
Termination	32	37

Table 2 Confusion matrix for recognizing initiation and termination behaviors using psychophysical features

Example / Classification	Initiation	Termination
Initiation	47	24
Termination	15	54

sensation is that the former embraces sensory experiences situated within the environment whereas the latter does not. Specifically, the psychophysical representation encodes the visual occlusion of the physical divider (described in Sect. 5.2) separating the 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 representation would mistakenly classify this as an adequate proxemic scenario (because it only encodes the distance), while the psychophysical representation would correctly classify this as an inadequate proxemic scenario (because the people are visually occluded). These “sensory interference” conditions are the hallmark of the continuous extension of the psychophysical representation proposed in our ongoing work [34].

7 Summary and Conclusions

In this work, we discussed a set of feature representations for analyzing human spatial behavior (proxemics) motivated by metrics used in the social sciences. Specifically, we considered individual, physical, and psychophysical factors that contribute to social spacing. We demonstrated the feasibility of autonomous real-time annotation of these prox-

emic features during a social interaction between two people and a humanoid robot in the presence of a visual obstruction (a physical barrier). We then used two different feature representations—physical and psychophysical—to train HMMs to recognize spatiotemporal behaviors that signify transitions into (initiation) and out of (termination) a social interaction. We demonstrated that the HMMs trained on psychophysical features, which encode the sensory experience of each interacting agent, outperform those trained on physical features, which only encode spatial relationships. These results suggest a more powerful representation of proxemic behavior with particular implications in autonomous socially interactive systems.

The models used by the system presented in this paper utilize heuristics based on empirical measures provided by the literature [15, 36, 44], resulting in a discretization of the parameter space. We are investigating a more continuous psychophysical representation learned from data (with a focus on voice loudness, visual, and kinesthetic factors) for the development of robust proxemic controllers for robots situated in complex interactions (e.g., with more than two agents, or with individuals with hearing or visual impairments) and environments (e.g., with loud noises, low light, or visual occlusions) [34].

The proxemic feature extraction and behavior recognition systems are part of the Social Behavior Library in the USC Interaction Lab ROS repository.¹³

Acknowledgements This work is supported in part by an NSF Graduate Research Fellowship, as well as ONR MURI N00014-09-1-1031 and NSF IIS-1208500, CNS-0709296, IIS-1117279, and IIS-0803565 grants. We thank Louis-Philippe Morency for his insights in integrating his head pose estimation system [37] and in the experimental design process, and Mark Bolas and Evan Suma for their assistance in using the PrimeSensor, and Edward Kaszubski for his help in integrating the proxemic feature extraction and behavior recognition systems into the Social Behavior Library.

References

1. Adams L, Zuckerman D (1991) The effect of lighting conditions on personal space requirements. *J Gen Psychol* 118(4):335–340
2. Aiello J (1987) Human spatial behavior. In: *Handbook of environmental psychology*, Chap 12. Wiley, New York
3. Aiello J, Aiello T (1974) The development of personal space: proxemic behavior of children 6 through 16. *Human Ecol.* 2(3):177–189
4. Aiello J, Thompson D, Brodzinsky D (1981) How funny is crowding anyway? Effects of group size, room size, and the introduction of humor. *Basic Appl Soc Psychol* 4(2):192–207
5. Argyle M, Dean J (1965) Eye-contact, distance, and affiliation. *Sociometry* 28:289–304
6. Bailenson J, Blascovich J, Beall A, Loomis J (2001) Equilibrium theory revisited: mutual gaze and personal space in virtual environments. *Presence* 10(6):583–598

¹³<https://code.google.com/p/usc-interaction-software.ros>.

7. Burgoon J, Stern L, Dillman L (1995) *Interpersonal adaptation: dyadic interaction patterns*. Cambridge University Press, New York
8. Cassell J, Sullivan J, Prevost S (2000) *Embodied conversational agents*. MIT Press, Cambridge
9. Dempster A, Laird N, Rubin D (1977) Maximum likelihood from incomplete data via the EM algorithm. *J R Stat Soc* 39:1–38
10. Deutsch RD (1977) Spatial structurings in everyday face-to-face behavior: a neurocybernetic model. *The Association for the Study of Man-Environment Relations*, Orangeburg
11. Evans G, Wener R (2007) Crowding and personal space invasion on the train: please don't make me sit in the middle. *J Environ Psychol* 27:90–94
12. Feil-Seifer D, Mataric M (2011) Automated detection and classification of positive vs negative robot interactions with children with autism using distance-based features. In: *HRI, Lausanne*, pp 323–330
13. Geden E, Begeman A (1981) Personal space preferences of hospitalized adults. *Res Nurs Health* 4:237–241
14. Hall ET (1959) *The silent language*. Doubleday Co, New York
15. Hall E (1963) A system for notation of proxemic behavior. *Am Anthropol* 65:1003–1026
16. Hall ET (1966) *The hidden dimension*. Doubleday Co, Chicago
17. Hall ET (1974) *Handbook for proxemic research*. American Anthropology Assn, Washington
18. Hayduk L, Mainprize S (1980) Personal space of the blind. *Soc Psychol Q* 43(2):216–223
19. Hediger H (1955) *Studies of the psychology and behaviour of captive animals in zoos and circuses*. Butterworths Scientific Publications, Stoneham
20. Huettnerach H, Eklundh K, Green A, Topp E (2006) Investigating spatial relationships in human-robot interaction. In: *IROS, Beijing*
21. Jan D, Traum DR (2007) Dynamic movement and positioning of embodied agents in multiparty conversations. In: *Proceedings of the 6th international joint conference on autonomous agents and multiagent systems, AAMAS '07*, pp 14:1–14:3
22. Jan D, Herrera D, Martinovski B, Novick D, Traum D (2007) A computational model of culture-specific conversational behavior. In: *Proceedings of the 7th international conference on intelligent virtual agents, IVA '07*, pp 45–56
23. Jones S, Aiello J (1973) Proxemic behavior of black and white first-, third-, and fifth-grade children. *J Pers Soc Psychol* 25(1):21–27
24. Jones S, Aiello J (1979) A test of the validity of projective and quasi-projective measures of interpersonal distance. *West J Speech Commun* 43:143–152
25. Jung J, Kanda T, Kim MS (2013) Guidelines for contextual motion design of a humanoid robot
26. Kendon A (1990) *Conducting interaction—patterns of behavior in focused encounters*. Cambridge University Press, New York
27. Kennedy D, Gläscher J, Tyszka J, Adolphs R (2009) Personal space regulation by the human amygdala. *Nat Neurosci* 12:1226–1227
28. Kristoffersson A, Severinson Eklundh K, Loutfi A (2013) Measuring the quality of interaction in mobile robotic telepresence: a pilot's perspective. *Int J Soc Robot* 5(1):89–101
29. Kuzuoka H, Suzuki Y, Yamashita J, Yamazaki K (2010) Reconfiguring spatial formation arrangement by robot body orientation. In: *HRI, Osaka*
30. Lawson B (2001) *Sociofugal and sociopetal space, the language of space*. Architectural Press, Oxford
31. Llobera J, Spanlang B, Ruffini G, Slater M (2010) Proxemics with multiple dynamic characters in an immersive virtual environment. *ACM Trans Appl Percept* 8(1):3:1–3:12
32. Low S, Lawrence-Zúñiga D (2003) *The anthropology of space and place: locating culture*. Blackwell Publishing, Oxford
33. Mead R, Mataric M (2011) An experimental design for studying proxemic behavior in human-robot interaction. Tech Rep CRES-11-001, USC Interaction Lab, Los Angeles
34. Mead R, Mataric MJ (2012) Space, speech, and gesture in human-robot interaction. In: *Proceedings of the 14th ACM international conference on multimodal interaction, ICMI '12*, Santa Monica, CA, pp 333–336
35. Mead R, Atrash A, Mataric MJ (2011) Recognition of spatial dynamics for predicting social interaction. In: *HRI, Lausanne, Switzerland*, pp 201–202
36. Mehrabian A (1972) *Nonverbal communication*. Aldine Transaction, Piscataway
37. Morency L, Whitehill J, Movellan J (2008) Generalized adaptive view-based appearance model: integrated framework for monocular head pose estimation. In: *8th IEEE international conference on automatic face gesture recognition (FG 2008)*, pp 1–8
38. Mumm J, Mutlu B (2011) Human-robot proxemics: physical and psychological distancing in human-robot interaction. In: *HRI, Lausanne*, pp 331–338
39. Oosterhout T, Visser A (2008) A visual method for robot proxemics measurements. In: *HRI workshop on metrics for human-robot interaction, Amsterdam*
40. Pelachaud C, Poggi I (2002) Multimodal embodied agents. *Knowl Eng Rev* 17(2):181–196
41. Price G, Dabbs J Jr. (1974) Sex, setting, and personal space: changes as children grow older. *Pers Soc Psychol Bull* 1:362–363
42. Rabiner LR (1990) A tutorial on hidden Markov models and selected applications in speech recognition. In: *Readings in speech recognition*, pp 267–296
43. Satake S, Kanda T, Glas DF, Imai M, Ishiguro H, Hagita N (2009) How to approach humans?: Strategies for social robots to initiate interaction. In: *HRI*, pp 109–116
44. Schegloff E (1998) Body torque. *Soc Res* 65(3):535–596
45. Schöne H (1984) *Spatial orientation: the spatial control of behavior in animals and man*. Princeton University Press, Princeton
46. Shotton J, Fitzgibbon A, Cook M, Sharp T, Finocchio M, Moore R, Kipman A, Blake A (2011) Real-time human pose recognition in parts from single depth images. In: *CVPR*
47. Sommer R (1967) Sociofugal space. *Am J Sociol* 72(6):654–660
48. Takayama L, Pantofaru C (2009) Influences on proxemic behaviors in human-robot interaction. In: *IROS, St. Louis*
49. Torta E, Cuijpers RH, Juola JF, van der Pol D (2011) Design of robust robotic proxemic behaviour. In: *Proceedings of the third international conference on social robotics, ICSR'11*, pp 21–30
50. Trautman P, Krause A (2010) Unfreezing the robot: navigation in dense, interacting crowds. In: *IROS, Taipei*
51. Vasquez D, Stein P, Rios-Martinez J, Escobedo A, Spalanzani A, Laugier C (2012) Human aware navigation for assistive robotics. In: *Proceedings of the thirteenth international symposium on experimental robotics, ISER'12*, Québec City, Canada
52. Walters M, Dautenhahn K, Boekhorst R, Koay K, Syrdal D, Nehaniv C (2009) An empirical framework for human-robot proxemics. In: *New frontiers in human-robot interaction, Edinburgh*

Ross Mead is a Computer Science PhD student, former NSF Graduate Research Fellow, and a fellow of the Body Engineering Los Angeles program (NSF GK-12) in the Interaction Lab at the University of Southern California (USC). Ross graduated with his Bachelor's degree in Computer Science from Southern Illinois University Edwardsville (SIUE) in 2007. At SIUE, his research efforts focused on interactions between groups of robots for organization, task allocation, and resource management. As the trend in robotic systems places them in social proximity of human users, his interests have evolved to consider the complexities of how robots and humans will interact. At USC, his research now focuses on the principled design and modeling of fundamental social behaviors—specifically, body language—to enable rich

autonomy in socially interactive robots targeted at supporting assistive and educational needs.

Amin Atrash is a Computer Science postdoctoral researcher in the Interaction Lab at the University of Southern California. He received his PhD from McGill University in 2011, his MS and BS from the Georgia Institute of Technology in 2003 and 1999, respectively, and worked as a research scientist at BBN technologies from 2003 to 2005. His current research focuses on the application of machine learning techniques in robotics, centered on learning and decision-making in human-robot interaction and multi-modal interfaces. His research has addressed the use of probabilistic models for data fusion, generation of social content, and recognition of multi-modal inputs.

Maja J. Matarić is professor and Chan Soon-Shiong chair in Computer Science, Neuroscience, and Pediatrics at the University of South-

ern California, founding director of the USC Center for Robotics and Embedded Systems (cres.usc.edu), co-director of the USC Robotics Research Lab (robotics.usc.edu), and Vice Dean for Research in the USC Viterbi School of Engineering. She received her PhD in Computer Science and Artificial Intelligence from MIT in 1994, MS in Computer Science from MIT in 1990, and BS in Computer Science from the University of Kansas in 1987. Her Interaction Lab's research into socially assistive robotics is aimed at endowing robots with the ability to help people through individual non-contact assistance in convalescence, rehabilitation, training, and education. Her research is currently developing robot-assisted therapies for children with autism spectrum disorders, stroke and traumatic brain injury survivors, and individuals with Alzheimer's Disease and other forms of dementia. Details about her research are found at <http://robotics.usc.edu/interaction/>.