

# A Preliminary Validation of Automated Feature Recognition for Proxemic Behavior Analysis

Ross Mead, Amin Atrash, and Maja J Matarić

Interaction Lab

Computer Science Department

Viterbi School of Engineering

University of Southern California

Los Angeles, CA 90089

{rossmead, atrash, mataric}@usc.edu

## Abstract

In this report, we present a preliminary evaluation of the feature extraction system for proxemic behavior analysis described in our previous work (Mead *et al.*, 2011a).

## 1 Introduction and Background

*Proxemics* is the study of the dynamic process by which people position themselves in face-to-face social encounters (Hall, 1959, 1963, 1966, 1974). This process is governed by sociocultural norms that, in effect, determine the overall sensory experience of each interacting participant (Hall, 1974). 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 (Deutsch, 1977; Mead *et al.*, 2011b; Schegloff 1998; Mehrabian, 1972). 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 (Oosterhout and Visser 2008; Jones and Aiello 1979). Recent developments in markerless motion capture — such as the Microsoft Kinect<sup>1</sup>, the Asus Xtion<sup>2</sup>, and the PrimeSensor<sup>3</sup> — have addressed the problem

of real-time human pose estimation, providing the means to revisit and more accurately model the subtle dynamics of human spatial interaction.

In this previous work, we presented a system that took advantage of these advancements to automate the analysis of proxemic behavior based on existing metrics (features) in the social sciences (Mead *et al.*, 2011a). In this report, we present a preliminary evaluation of that system.

## 2 System Implementation and Evaluation

The feature extraction system described in (Mead *et al.*, 2011a) can be implemented using any human motion capture technique. We utilized the PrimeSensor structured light range sensor and the OpenNI<sup>4</sup> 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 (Schegloff, 1998), physical (Mehrabian, 1972), and psychophysical (Hall, 1963) features, discussed below. Each feature set is described in detail in (Mead *et al.*, 2011a).

---

<sup>1</sup> <http://www.xbox.com/kinect>

<sup>2</sup> [http://www.asus.com/Multimedia/Motion\\_Sensor](http://www.asus.com/Multimedia/Motion_Sensor)

<sup>3</sup> <http://www.primesense.com>

---

<sup>4</sup> <http://openni.org>

## 2.1 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 technique<sup>5</sup> 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):  $E(d) = k \times d^2$ , where  $k = 0.0055$ . Our procedure and results are consistent with that reported in the ROS Kinect accuracy test<sup>6</sup>, as well as that 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  $k$ .

## 2.2 Individual Features

Shotton *et al.* (2011) 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.

## 2.3 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 (Shotton *et al.*, 2011), 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 (Shotton *et al.*, 2011).

## 2.4 Psychophysical Features

Each feature annotation in Hall's (1963) psychophysical representation was developed based upon values from literature on the human sensory system (Hall, 1966). It is beyond the scope of this work to evaluate whether or not a participant actually experiences the stimulus in the way specified by a particular feature interval. For example, we do not measure the radiant heat or odor transmitted by one individual and the corresponding intensity at the sensory organ of the receiving individual. Rather, this work seeks to validate the annotations of these psychophysical features as a function of their respective distance and orientation intervals based on the sensor error models developed throughout this section; thus, the following is a theoretical model of the error associated with each psychophysical feature code.

Intervals for each feature code were evaluated at 1, 2, 3, and 4 meters from the sensor; the sensor was orthogonal to a line passing through the hip poses of two standing individuals. The results of these estimates are illustrated in Figures 1-8. At some ranges from the sensor, a feature interval would place the individuals out of the field-of-view; these values are omitted<sup>7</sup>.

The evaluation of feature intervals for the *distance*, *olfaction*, *touch*, *thermal*, and *voice loudness codes* demonstrates the effect of interval size on feature classification (Figs. 1-5, respectively). Specifically, the larger and/or more common feature intervals are extracted with higher precision, suggesting their potential utility for annotating proxemics in typical social encounters.

The *SFP axis code* contains feature intervals of uniform size (20°). 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

---

<sup>5</sup> [http://www.pointclouds.org/documentation/tutorials/planar\\_segmentation.php](http://www.pointclouds.org/documentation/tutorials/planar_segmentation.php)

<sup>6</sup> [http://www.ros.org/wiki/openni\\_kinect/kinect\\_accuracy](http://www.ros.org/wiki/openni_kinect/kinect_accuracy)

---

<sup>7</sup> This 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.

individual with respect to 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. 6).

Eye gaze was unavailable for true *visual code* extraction. Instead, the *visual code* is estimated based on the shoulder orientation, resulting in a coarse estimates of narrow feature intervals (foveal and macular; Fig. 7).

For evaluation of the *kinesthetic code*, we used interval distance estimates based on average human limb lengths (Fig. 8; Hall, 1966). In practice, performance is expected to be lower. The feature intervals are variable, as 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.

### 3 Conclusions and Future Work

In this report, we presented a preliminary evaluation of the feature extraction system for proxemic behavior analysis described in our previous work (Mead *et al.*, 2011a). We utilized a markerless motion capture system (the PrimeSensor) to recognize these features 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. Our results demonstrate the feasibility of autonomous real-time annotation of these proxemic features.

Currently, the models in our automated system utilize heuristics based on metrics provided by the literature (Mead *et al.*, 2011a); model parameters have been determined from empirical evidence as observed by these researchers. Consequently, this approach results in a discretization of the parameter space. Contemporary machine learning techniques could be used to automatically determine these parameters from data.

### Acknowledgments

This work is supported in part by the National Science Foundation under the Graduate Research Fellowship Program, as well as ONR MURI N00014-09-1-1031 and NSF CNS-0709296 grants.

We thank Mark Bolas and Evan Suma for their assistance in using the PrimeSensor.

### References

- Albert Mehrabian. 1972. *Nonverbal Communication*. Aldine Transaction, Piscataway, New Jersey.
- Edward T. Hall. 1959. *The Silent Language*. Doubleday Co., New York, New York.
- Edward T. Hall. 1963. A System for Notation of Proxemic Behavior. *American Anthropologist*, 65:1003-1026.
- Edward T. Hall. 1966. *The Hidden Dimension*. Doubleday Co., New York, New York.
- Edward T. Hall. 1974. *Handbook for Proxemic Research*. American Anthropology Assn., Washington, District of Columbia.
- Emanuel A. Schegloff. 1998. Body Torque. *Social Research*, 65(3):535-596.
- Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard Moore, Alex Kipman, Andrew Blake. 2011. Real-time Human Pose Recognition in Parts from Single Depth Images. 24th IEEE Conference on Computer Vision and Pattern Recognition, Colorado Springs, Colorado.
- Robert D. Deutsch. 1977. *Spatial Structurings in Everyday Face-to-face Behavior*. The Association for the Study of Man-Environment Relations, Orangeburg, New York.
- Ross Mead, Amin Atrash, and Maja J Matarić. 2011a. Proxemic Feature Recognition for Interactive Robots: Automating Metrics from the Social Sciences. *International Conference on Social Robotics*, 52-61.
- Ross Mead, Amin Atrash, and Maja J Matarić. 2011b. Recognition of Spatial Dynamics for Predicting Social Interaction. 6th ACM/IEEE International Conference on Human-Robot Interaction, 201-202.
- Ross Mead and Maja J Matarić. 2011. An Experimental Design for Studying Proxemic Behavior in Human-Robot Interaction. Technical Report, CRES-11-001.
- Stanley E. Jones and John R. Aiello. 1979. A Test of the Validity of Projective and Quasi-Projective Measures of Interpersonal Distance. *The Western Journal of Speech Communication*, 43:143-152.
- Tim van Oosterhout and Arnoud Visser. 2008. A Visual Method for Robot Proxemics Measurements. 3rd ACM/IEEE International Conference on Human-Robot Interaction Workshop on Metrics for Human-Robot Interaction, Technical Report 471, Amsterdam, Netherlands.

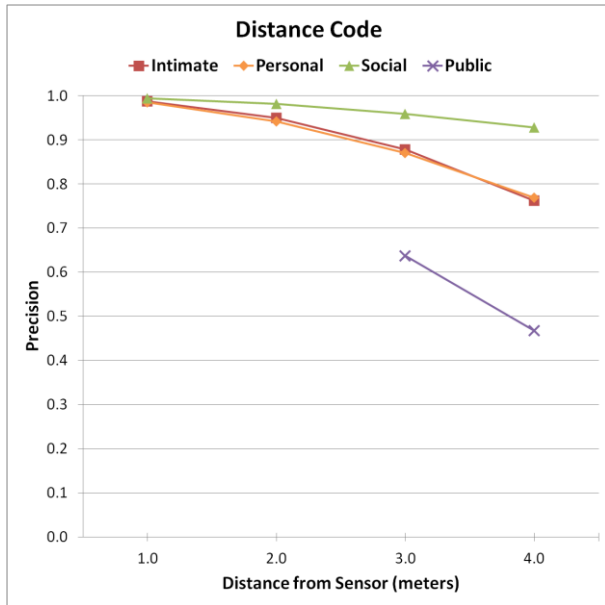


Figure 1: Error model of psychophysical distance code.

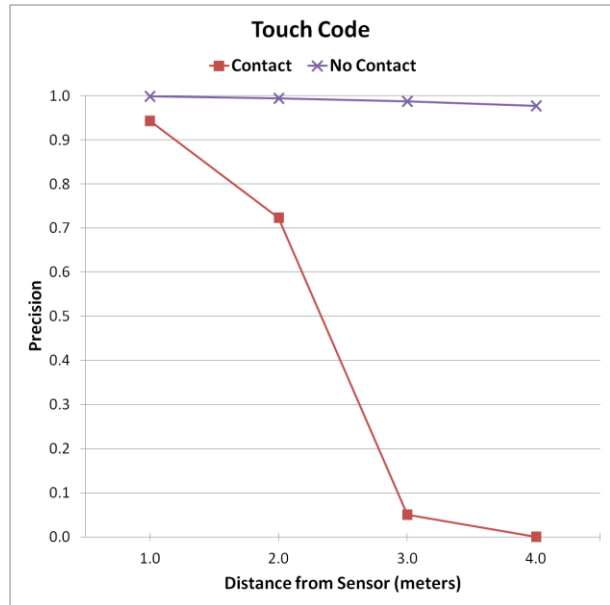


Figure 3: Error model of psychophysical touch code.

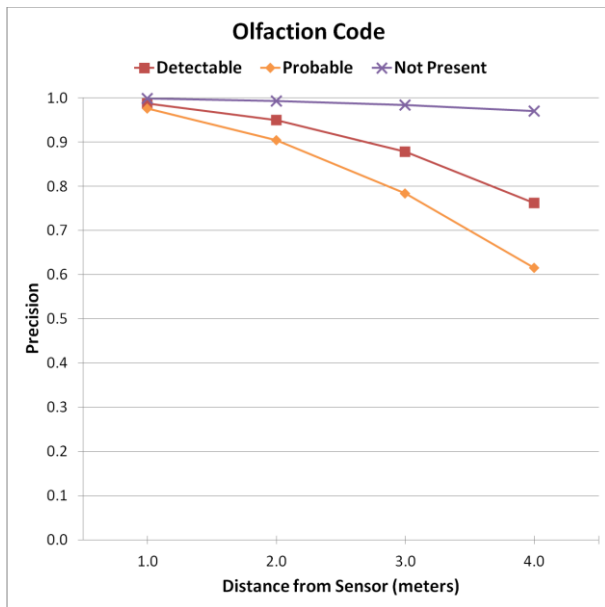


Figure 2: Error model of psychophysical olfaction code.

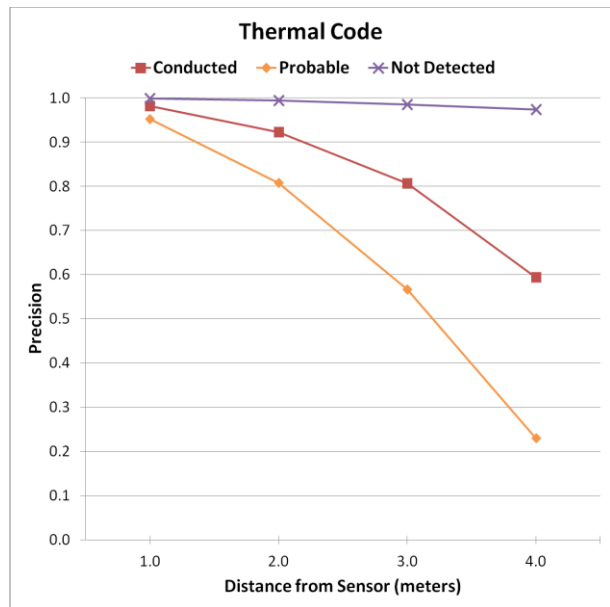


Figure 4: Error model of psychophysical thermal code.

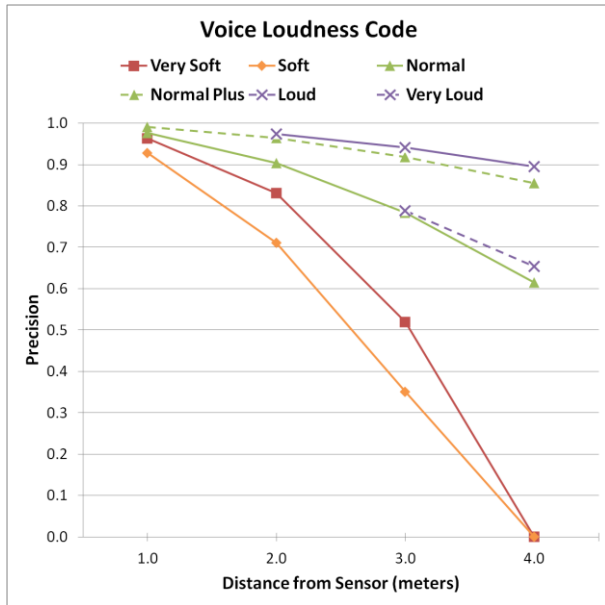


Figure 5: Error model of psychophysical voice loudness code.

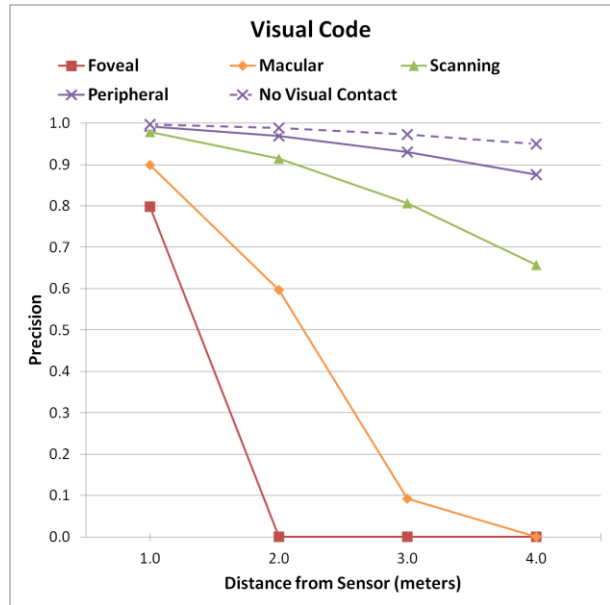


Figure 7: Error model of psychophysical visual code.

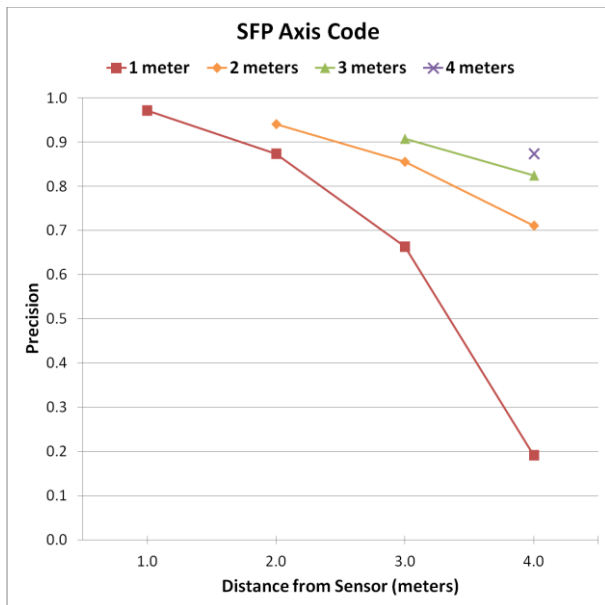


Figure 6: Error model of psychophysical sociofugal-sociopetal (SFP) axis code.

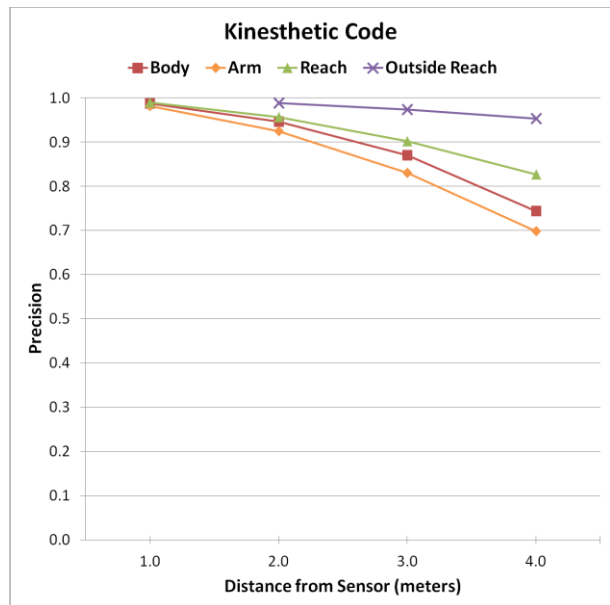


Figure 8: Error model of psychophysical kinesthetic code.