

Mediated Social Interpersonal Communication:
Evidence-based Understanding of Multimedia Solutions for Enriching Social Situational

Awareness

by

Sreekar Krishna

A Dissertation Presented in Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy

Approved March 2011 by the
Graduate Supervisory Committee:

Sethuraman Panchanathan, Chair

John Black Jr.

Baoxin Li

Gang Qian

Michelle Shiota

ARIZONA STATE UNIVERSITY

May 2011

ABSTRACT

Social situational awareness, or the attentiveness to one's social surroundings, including the people, their interactions and their behaviors is a complex sensory-cognitive-motor task that requires one to be engaged thoroughly in understanding their social interactions. These interactions are formed out of the elements of human interpersonal communication including both verbal and non-verbal cues. While the verbal cues are instructive and delivered through speech, the non-verbal cues are mostly interpretive and requires the full attention of the participants to understand, comprehend and respond to them appropriately. Unfortunately certain situations are not conducive for a person to have complete access to their social surroundings, especially the non-verbal cues. For example, a person is who is blind or visually impaired may find that the non-verbal cues like smiling, head nod, eye contact, body gestures and facial expressions of their interaction partners are not accessible due to their sensory deprivation. The same could be said of people who are remotely engaged in a conversation and physically separated to have a visual access to one's body and facial mannerisms. This dissertation describes novel multimedia technologies to aid situations where it is necessary to mediate social situational information between interacting participants.

As an example of the proposed system, an evidence-based model for understanding the accessibility problem faced by people who are blind or visually impaired is described in detail. From the derived model, a sleuth of sensing and delivery technologies that use state-of-the-art computer vision algorithms in combination with novel haptic interfaces are developed towards a) *A Dyadic Interaction Assistant*, capable of helping individuals who are blind to access important head and face based non-verbal communicative cues during one-on-one dyadic interactions, and b) *A Group Interaction Assistant*, capable of provide situational awareness about the interaction partners and their dynamics to a user who is blind, while also providing important social feedback about their own body mannerisms. The goal is to increase the effective social situational information that one has access to, with the conjuncture that a good awareness of one's social surroundings gives them the ability to understand and empathize with their interaction partners better. Extending the work from an important social interaction assistive technology, the need for enriched social situational awareness

is everyday professional situations are also discussed, including, a) enriched remote interactions between physically separated interaction partners, and b) enriched communication between medical professionals during critical care procedures, towards enhanced patient safety.

In the concluding remarks, this dissertation engages the readers into a science and technology policy discussion on the potential effect of a new technology like the social interaction assistant on the society. Discussing along the policy lines, social disability is highlighted as an important area that requires special attention from researchers and policy makers. Given that the proposed technology relies on wearable inconspicuous cameras, the discussion of privacy policies is extended to encompass newly evolving interpersonal interaction recorders, like the one presented in this dissertation.

To my grandfather.

ACKNOWLEDGEMENTS

I would like to express my deepest gratitude to my committee chair, Dr. Sethuraman Panchathan, for being an excellent mentor, guide and role model towards developing not only my intellectual skills, but also the larger skill set of how to achieve success in life and career. I sincerely thank him for showing at most patience in allowing me to search for my passion and continue along lines of research which are otherwise considered closed for typical Electrical Engineering students.

I would also like to thank Dr. John Black for his invaluable advice and suggestions, both during my research and during the writing of this dissertation. John has made my research experience at Arizona State University a pleasurable one.

I sincerely acknowledge the role Dr. Terri Hedgpeth played in shaping my research career for the past 5 years. She was in essence the guiding light for my research and all the designs I conceptualized.

I would like to convey my sincere thanks to Dr. Michelle Shiota for accepting the challenge of mentoring a computer science student in the fascinating and deep areas of facial expression recognition and human social interactions. I want to extend my thanks to all her lab members who patiently helped me in my research, especially Wan Yeung for enrolling me in her FACS training class and taking the time to explain concepts that are otherwise unfamiliar to non-psychology students.

I thank Dr. Gang Qian and Dr. Baoxin Li for serving on my dissertation committee, while providing guidance in my research.

I would like to thank all my colleagues and fellow researchers, at the Center for Cognitive Ubiquitous Computing for helping me, and for providing valuable advice.

TABLE OF CONTENTS

	Page
LIST OF TABLES	xvi
LIST OF FIGURES	xviii
CHAPTER	1
1 SITUATIONAL AWARENESS IN EVERYDAY SOCIAL INTERACTIONS	1
1.1 Components of Social Interactions	2
1.1.1 Non-verbal communication cues	3
1.1.1.1 Social Sight and Social Hearing	4
1.1.1.2 Social Touch	5
1.2 Social Situational Awareness	5
1.2.1 Social Situational Awareness in Everyday Social Interactions	7
1.2.1.1 SSA in Dyadic Interactions	7
1.2.1.2 SSA in Group Interactions	7
1.2.2 Learning Social Awareness	8
1.3 Factors that are important ro successful non-verbal communication	10
1.3.1 Factors related to the Communication Environment	10
1.3.2 Factors related to the Physical Characteristics of the communicators	11
1.3.3 Factors related to the Behaviors of the Communicator	12
1.3.3.1 Gesture	12
1.3.3.2 Posture	13
1.3.3.3 Touch	13
1.3.3.4 Face/Head	14
1.3.3.5 Eye	16
1.4 Facilitating Social Interactions through the Enrichment of Social Situational Awareness	16
1.5 Organization of the Dissertation	17
2 NEED FOR ENRICHING SOCIAL SITUATIONAL AWARENESS	19
2.1 Disability Induced Social Signal Attrition	19

Chapter	Page
2.1.1 Visual Impairment - a hinderance to smooth social interactions . . .	20
2.1.1.1 Inability to learn social skills due to the lack of vision: . . .	21
2.1.1.2 Lack of visual reinforcement feedback on one's manner- isms:	22
2.2 Social Signal Attrition during Remote Interpersonal Interactions	23
2.3 Social Signal Attrition in Medical Teams	26
2.3.1 Importance factors affecting team performance	29
2.3.1.1 Group Dynamics	29
2.3.1.2 Leadership	30
2.3.2 The emerging science of medical team social assessment	31
2.3.2.1 The Mayo Clinic Multidisciplinary Simulation Center: . .	32
2.3.2.2 Challenges in Social Situational Assessment and Training	33
2.4 The Research Focus of this Dissertation	35
2.4.1 The Handshake Example: An Example of Evidence-based Under- standing of Social Situational Awareness	36
3 ASSISTIVE MEDIATION TECHNOLOGY FOR INDIVIDUALS WHO ARE VISUALLY IMPAIRED	39
3.1 Automatic Detection of Non-verbal Cues and Observations	39
3.1.1 Design principles for social assistive and rehabilitative devices . . .	42
3.2 Related Work in Computer Vision Research towards Sensing Factors that Contribute to the Overall Non-verbal Communication Picture	43
3.3 Requirements Analysis for a Social Assistive Technology: Evidence Ag- gregation	45
3.3.1 Results from the Online Survey	47
3.3.1.1 Average Response	47
3.3.1.2 Response on Individual Questions	47
3.3.1.3 Response Ratio	48
3.3.1.4 Rank Average Importance Map for Various Non-verbal Cues	49

Chapter	Page
3.4 Evidence-based Model for the Proposed Social Interaction Assistant	51
3.5 System Architecture for a Social Interaction Assistant: An Evidence-based Assessment of Requirements	52
3.6 Organization of the Later Chapters	53
4 ENRICHING SOCIAL FEEDBACK: SENSING STEREOTYPIC BEHAVIORS	55
4.1 Stereotypic Body Behaviors	55
4.2 Focus of the chapter	57
4.3 Background and Related Work	58
4.3.1 Foundations for social rehabilitation of behavioral stereotypes . . .	58
4.3.1.1 <i>Intervention</i>	58
4.3.1.2 <i>Self Monitoring</i>	58
4.3.2 Need for Assistive or Rehabilitative Technology	60
4.3.2.1 Past research into building assistive technology to detect body rocking	60
4.4 Methodology	62
4.4.1 Motion Sensors - Design choice along the “Acceptance” Dimension	64
4.4.2 Extracting Body Rock Information from Motion Sensor Data - De- sign choice along the “Motivation” Dimension	66
4.4.2.1 Features:	67
4.4.2.2 Learning Algorithm:	71
4.5 Data Collection	74
4.5.1 Controlled Data Collection:	74
4.5.1.1 Routine A: Rocking data	75
4.5.1.2 Routine B: Non-rocking data	75
4.5.1.3 Routine C: Test data	75
4.5.2 Uncontrolled Data Collection:	75
4.6 Experiments	76
4.7 Results	77
4.8 Discussion of Results	84

Chapter	Page
4.8.1 Packet Length, and Detection Efficiency	85
4.8.2 Generalization Capabilities	86
5 SENSING FACIAL EXPRESSIONS IN DYADIC INTERACTIONS	89
5.1 Design Considerations	90
5.2 Dyadic Interaction Assistant - Proposed Solution	91
5.3 Facial Expression Recognition - State of the art	92
5.4 Design Considerations for Dyadic Interaction Assistant	95
5.5 Temporal Exemplar-based Bayesian Network (TEBN) for Facial Expression & Gesture Sensing	96
5.5.1 The Observation Layer	98
5.5.2 The Exemplar Layer	100
5.5.2.1 Computing $P(X(t)/L_{ij}(t))$: Representing the test data in terms of existing examples	101
5.5.2.2 Computing $P(L_i(t)/Y_i)$: Determining the prior probability of chosen examples w.r.t the complete training data	102
5.5.3 The Prior Knowledge Layer	103
5.5.3.1 Computation of $P(Y_i/H(t))$: Temporal propagation of expression probabilities	104
5.6 Discussion of Design Considerations: TEBN perspectives	105
5.7 Experiments	106
5.7.1 Data	106
6 DELIVERING DYADIC INTERACTION CUES THROUGH VIBROTACTILE STIMULATIONS	108
6.1 Related Work in Haptic Interpersonal Communication Interfaces	109
6.2 Proposed Visuo-Haptic Sensory Substitution Device	114
6.3 The Vibrotactile Glove	118
6.4 Haptic Cueing to Test Localization and Spatio-Temporal Mapping	120
6.4.1 Localization	120
6.4.2 Spatio-Temporal Cueing	121

Chapter	Page
6.4.2.1	Group 1 - The visual emoticon motivated haptic icons: . . . 121
6.4.2.2	Group 2 - The auxiliary haptic icons: 123
6.5	Research Hypothesis 123
6.5.1	Localization 123
6.5.1.1	Hypothesis 1: 123
6.5.1.2	Hypothesis 2: 124
6.5.1.3	Hypothesis 3: 124
6.5.2	Spatio-Temporal Cueing 124
6.5.2.1	Hypothesis 4: 124
6.5.2.2	Hypothesis 5: 124
6.6	Experiments and Analysis Methodology 125
6.6.1	Participants 125
6.6.2	Procedure 125
6.6.3	Analysis 126
6.6.3.1	Recognition Accuracies 127
6.6.3.2	ANOVA 127
6.6.3.3	Tuckey HSD 127
6.7	Results of the Experiments 128
6.7.1	Localization Experiments 129
6.7.1.1	Phalange Level Localization 129
6.7.1.2	Finger Level Localization 130
6.7.1.3	Phalange Position Localization 131
6.7.2	Spatio-Temporal Experiments 132
6.7.2.1	Individual Spatio-temporal Patterns 132
6.7.2.2	Comparison of Spatio-temporal Cueing groups 133
6.7.2.3	Time for Recognition: 135
6.8	Conveying Facial Expressions through Dyadic Social Situational Assistant . 136
7	BIOMETRICS IN SOCIAL CONTEXT: IDENTIFYING INTERACTION PART- NERS 138

Chapter	Page
7.0.1 Employing face recognition to facilitate social interactions	139
7.1 Face Recognition in Humans	140
7.2 Our Approach to Face Recognition	143
7.3 Feature Extractors	143
7.3.1 Gabor Features	144
7.3.1.1 Use of Gabor Filters in Face Recognition	144
7.3.1.2 Gabor Filters	146
7.3.1.3 Gaussian Function	147
7.3.1.4 Sinusoid	147
7.4 The Learning Algorithm	149
7.4.1 Genetic Algorithms	150
7.4.1.1 Use of Genetic Algorithms in Face Recognition	151
7.4.1.2 The Chromosome	153
7.4.1.3 Creation of the first generation	154
7.4.1.4 Creation of the newer generations	156
7.4.1.5 <i>Crossover</i>	157
7.4.1.6 <i>Mutation</i>	157
7.5 Methodology	158
7.5.1 The FacePix (30) Database	159
7.5.2 The Gabor Features	161
7.5.3 The Genetic Algorithm	162
7.5.3.1 The Fitness Function	163
7.6 Results	167
7.6.1 Discussion of Results	168
7.6.2 Person-specific feature extraction	169
7.7 Face Recognition in Social Situational Awareness Context	170
8 ENRICHING GROUP INTERACTIONS: SENSING PROXEMICS AND THE SOCIAL SPACE	172
8.1 Accurate Face Detection	173

Chapter	Page
8.2 Related Work in Accurate Face Detection	176
8.3 Proposed Framework	176
8.3.1 Module 1: Human Skin Tone Detector with Dynamic Background Modeler	177
8.3.1.1 <i>a-priori</i> Bi-modal Gaussian Mixture Model for Human Skin Classification	177
8.3.1.2 Dynamically Learnt Multi-modal Gaussian Model for Back- ground Pixel Classification	178
8.3.1.3 Skin and Background Classification using the learnt Multi- modal Gaussian Models	179
8.3.2 Module 2: Evidence-Aggregating Human Face Silhouette Random Field Modeler	180
8.3.2.1 Random Field (RF) Models	180
8.3.2.2 Pre-processing	181
8.3.2.3 The Neighborhood System	182
8.3.2.4 Local Conditional Probability Density (LCPD)	183
8.3.2.5 Human Face Pose	185
8.3.3 Combining Evidence	185
8.3.3.1 Dempster-Shafer Theory of Evidence (DST)	186
8.3.4 Coarse Pose estimation	187
8.4 Testing the Abilities of the Face Detector	188
8.5 Results	188
8.6 Discussion of Results	190
9 SENSING DYNAMICS OF THE SOCIAL SCENE	192
9.1 Challenges in Person Localization from a wearable camera platform	193
9.1.1 Background Properties	193
9.1.2 Object Properties	193
9.1.3 Object/Camera Motion	194
9.1.4 Other Important Factors Affecting Effective Person Tracking	195

Chapter	Page
9.2 Related Computer Vision Work in Person Localization and Tracking	196
9.2.1 Detection Algorithms	196
9.2.2 Tracking Algorithms	197
9.3 Conceptual Framework	199
9.4 Structured Mode Searching Particle Filter	200
9.4.1 Step 1: Particle Filtering Step	200
9.4.2 Step 2: Structured Search	204
9.4.3 Chamfer Matching in Structured Search	207
9.5 Experiments and Datasets	209
9.5.1 Datasets	209
9.5.2 Evaluation Metrics	209
9.6 Results	211
10 COMMUNICATING SOCIAL SCENE DYNAMICS	216
10.1 Proposed Framework of Social Scene Structure Delivery	216
10.2 Related Work in Haptic Vibrotactile Technology for Information Delivery .	218
10.3 Design Requirements	222
10.4 Implementation	224
10.4.1 Form Factor	224
10.4.2 System Architecture	225
10.5 Hardware Design	226
10.5.1 Control Box	226
10.5.1.1 Main Controller	227
10.5.1.2 Bus Communication	227
10.5.1.3 Power Supply	227
10.5.1.4 Wireless Hardware	227
10.6 Tactor Modules	228
10.7 Software Design	228
10.7.1 Firmware	228
10.7.1.1 Main Controller Firmware	229

Chapter	Page
10.7.1.2 The Tactor Controller	230
10.8 User Interface	234
10.9 Experiments	235
10.9.1 Experiment 1: Localization of Vibrotactile Cues	235
10.9.1.1 Apparatus:	235
10.9.1.2 Procedure:	236
10.9.1.3 Results:	236
10.9.1.4 Discussion:	237
10.9.2 Experiment 2: Signal Duration as Cue for Distance	238
10.9.2.1 Subjects:	238
10.9.2.2 Apparatus:	238
10.9.2.3 Procedure:	238
10.9.2.4 Results:	239
10.9.2.5 Discussion:	239
10.9.3 Experiment 3: Vibrotactile Rhythm as Cue for Distance	241
10.9.3.1 Tactile Rhythm Design	241
10.9.3.2 Experiment	242
10.10 Other Applications for the proposed technology	246
10.10.1 Navigation and Spatial Orientation	246
10.10.2 Interpersonal Social Communication	247
10.10.3 Generic Information Communication	247
10.10.4 Case Study: Waist-worn Vibrotactile Display for Pedagogical Ap- plication for Choreographed Dance	248
10.10.4.1 Related Work in the Use of Vibrotactile Cues for Teach- ing Dance	248
10.10.4.2 Subjects	249
10.10.4.3 Procedure	250
10.10.4.4 Aim	252
10.10.4.5 Results	255

10.11 Group Interaction Dynamics Through Haptic Belt	259
11 THE SOCIAL INTERACTION ASSISTANT & ITS ROLE IN SOCIETY	261
11.1 Role of Social Assistance in the Society: A Policy Discussion	263
11.1.1 Social Disability: The hidden barrier to professional growth in the disabled population	265
11.1.2 Effect on Social & Emotional Intelligence due to Disability	267
11.1.3 Psychological Breakdown related to Social Skills	269
11.1.4 Societal Metrics of Social Disability	271
11.1.4.1 Comparison of Economic Status	271
11.1.4.2 Comparison of Personal Lives	271
11.2 Wearable Cameras: Ethics & Privacy	274
12 Conclusions & Future Work	280
REFERENCES	282
APPENDIX	314
A ALGORITHM FOR ESTIMATING RANK AVERAGE OF GROUPS	315
A.0.1 Procedure	316
B CONVEX OPTIMIZATION USING NEWTON'S METHOD - ENTROPY MAX- IMIZATION UNDER LINEAR CONSTRAINTS	317
B.1 Entropy Maximization under Linear Constraints Problem:	318
B.2 Newton's method	321
B.2.1 Newton's Method for Estimating Weights w_{ij}	322
C AMERICAN COMMUNITY SURVEY FORM - SAMPLE PAGES FROM 2008 SURVEY FORM	324

LIST OF TABLES

Table	Page
1.1 Eight factors of the environment that can affect interpersonal communication. . .	11
1.2 The physical characteristics of a communicator that can affect interpersonal communications.	12
1.3 FACS communicative actions on the human face	15
1.4 The role of the eyes in human interpersonal communication.	17
1.5 Organization of the dissertation.	18
2.1 Survey on the challenges of remote interaction.	27
3.1 Related Work in Automatic Detection of Interaction Environment Cues	43
3.2 Related Work in Automatic Detection of 1) Physical Characteristics of the Communicator, and 2)Behavior of the Communicator	44
3.3 Average Likert Score on the 8 statements obtained through an Online Survey. .	48
3.4 Organization of the dissertation.	54
4.1 Features for Body Rock Detection: Group 1	69
4.2 Features for Body Rock Detection: Group 2	70
4.3 Experiments with naturalistic data	79
5.1 State-of-the-art facial expression recognition algorithms and their performance. Exp: Spontaneous(S) / Posed expression (P); Per: Person Dependent(P) / Independent (I); Class: Number of classes; Sub: Number of subjects; Type: Data Type: Video (V) / Image(I); % Acc: Percentage Accuracy; ?: missing entry. . .	94
6.1 Haptic interpersonal interaction enrichment devices	110
6.2 Haptic interpersonal interaction enrichment devices contd.	111
6.3 Haptic interpersonal interaction enrichment devices contd.	112
6.4 Haptic interpersonal interaction enrichment devices contd.	113
6.5 Related work in the development of haptic gloves for information delivery . . .	115
6.6 Related work in the development of haptic gloves for information delivery contd.	116
6.7 Related work in the development of haptic gloves for information delivery contd.	117
8.1 Face detection validation results on FERET database.	189

Table	Page
8.2 Face detection validation results on the in-house face database.	189
10.1 Design Requirements for Vibrotactile Belts	222
10.2 Foot Steps Involved in the Choreographed Dance Movements.	252

LIST OF FIGURES

Figure	Page
1.1 Relative importance of a) verbal vs. non-verbal cues, b) four channels of non-verbal cues, and c) visual vs. audio encoding and decoding of bilateral human interpersonal communicative cues. Based on the meta-analysis presented in [1].	3
1.2 Relative communicative information plotted against its leakiness. Speech forms the verbal channel. Face, body and voice form the non-verbal communication channels.	4
1.3 Social situational awareness in human social communications	6
1.4 Social learning systems with continuous learning feedback loop.	9
1.5 The facial muscles responsible for facial expressions and gestures.	16
2.1 TeamSTEPPS: Team Strategies and Tools to Enhance Performance and Patient Safety	28
2.2 Multi goal model of self regulation for effective team leadership.	30
2.3 Evidence-based Model for Social Situation Awareness to promote handshake non-verbal cueing in the visually impaired and blind population.	38
3.1 Histogram of Responses grouped by Questions	49
3.2 Response Ratio	50
3.3 Rank average of the 8 questions	51
3.4 Social Interaction Assistant System Architecture.	52
4.1 Training and testing phases of a typical learning framework found in literature.	63
4.2 The proposed hardware for use in the detection of body rocking stereotypic behavior. The accelerometer, in comparison with a US quarter, is shown in the inset. The three axes marked in the image shows the orientation of the accelerometer as it is placed on the head.	65
4.3 Data stream for the tri-axial accelerometer. The three streams correspond to the three axes. The figure shows non-rocking events followed by rocking and then followed by non-rocking.	66

Figure	Page
4.4 Packet length to recognition rate comparison under the classic AdaBoost framework.	78
4.5 Packet length to recognition rate comparison under the Modest AdaBoost framework.	78
4.6 Piecewise performance analysis of the classic AdaBoost classifier framework; (a) Recognition rates under use of individual feature sets; (b) The Receiver Operating Characteristics (ROC) under the use of individual feature sets; (c) Area under the curve (AUC) for each feature set as estimated from the ROC; (d) The number of simple classifiers used by the aggregated AdaBoost classifier. Each set and each feature representation in the classifier pool are separately marked. In all the graphs Set 1 through 5 are as explained by Tables 4.1 and 4.2. Set 6 represents a set containing all 14 features from Tables 4.1 and 4.2.	81
4.7 Piecewise performance analysis of the Modest AdaBoost framework; (a) Recognition rates under use of individual feature sets; (b) The Receiver Operating Characteristics (ROC) under the use of individual feature sets; (c) Area under the curve (AUC) for each feature set as estimated from the ROC; (d) The number of simple classifiers used by the aggregated AdaBoost classifier; Each set and each feature representation in the classifier pool are separately marked. In all the graphs Set 1 through 5 are as explained by Tables 4.1 and 4.2. Set 6 represents a set containing all 14 features from Tables 4.1 and 4.2.	83
5.1 Typical use of the dyadic interaction assistance scenario, a third person perspective on the use case scenario.	91
5.2 Face tracker with an auto-focus camera and a micro pan-tilt mechanism.	92
5.3 Typical use of the dyadic interaction assistance scenario, a third person perspective on the use case scenario.	93
5.4 Example Bayesian Network.	97
5.5 Temporal Exemplar-based Bayesian Network for facial expression recognition.	99
5.6 36 Facial fiducial points tracked with FaceAPI software. Both x and y coordinates from all 36 points are used for facial expression recognition.	99

Figure	Page
5.7 Deriving the exemplar layer of the TEBN based on every test point $X_t(t)$	101
5.8 Comparison of Fear and Surprise facial expressions.	104
6.1 Somatosensory homonculus as mapped through magnetoencephalography [2].	114
6.2 The Vibrotactile Glove.	118
6.3 Localization and spatio-temporal cueing software used for the vibrotactile glove.	119
6.4 Phalange naming convention and grouping based on the anatomical distances.	120
6.5 Mapping of Group 1 and Group 2 haptic expression icons to the central three fingers (9 Phalanges) of the vibrotactile glove. In the expression mapping chart, Columns 1 to 3 represent the expression. Column 4 shows the spatial mapping of vibrations. Column 5 shows the temporal mapping of the vibrations.	122
6.6 (a) Recognition Accuracies; (b) ANOVA; (c) HSD	129
6.7 (a) Recognition Accuracies; (b) ANOVA; (c) HSD	130
6.8 (a) Recognition Accuracies; (b) ANOVA; (c) HSD	131
6.9 (a) Recognition Accuracies; (b) ANOVA; (c) HSD	132
6.10 (a) Recognition Accuracies; (b) ANOVA; (c) HSD	133
6.11 Confusion Matrix across the 12 participants. The rows are the stimulation and the columns are the responses of the participants. Each row adds to 100% (rounding error of 1%).	135
6.12 Average recognition rate and response time for the subject who is blind, for over 70 trails.	136
6.13 Average response time for all 12 participants. Four important results are shown above, 1) Avg. correct response time per expression (Cyan), 2) Avg. incorrect response time per expression (Red), 3) Avg. correct response time for Group 1 (Blue), and 4) Avg. correct response time for Group 2 (Magenta).	137
7.1 (a) 3D representation of a Gaussian mask; $\sigma_x = 10$, $\sigma_y = 15$ and $\theta = 0$ (b)Image of the Gaussian mask $\sigma_x = 10$, $\sigma_y = 15$ and $\theta = 0$	147
7.2 (a)3D representation of a Sinusoid $S_{\omega,\theta}$ (b)Image representation of the real part of the complex Sinusoid $\Re\{S_{\omega,\theta}\}$ (c)Image representation of the imaginary part of complex Sinusoid $\Im\{S_{\omega,\theta}\}$	148

Figure	Page
7.3 (a)3D representation of a Gabor filter $\Psi_{\omega,\theta}$	
(b)Image representation of the real part of Gabor filter $\Re\{\Psi_{\omega,\theta}\}$	
(c)Image representation of the imaginary part of Gabor filter $\Im\{\Psi_{\omega,\theta}\}$	149
7.4 A typical chromosome used in the proposed method.	153
7.5 Stages in the creation of the first generation of parents	154
7.6 Deriving newer parents from the current generation	155
7.7 Typical crossing of two parents to create an offspring	157
7.8 Mutation of a newly created offspring	158
7.9 The data capture setup for FacePix(30)	159
7.10 Sample face images with varying pose and illumination from the FacePix(30) database	160
7.11 Sample frontal images of one person from the FacePix(30) Database	161
7.12 A face image marked with 5 locations where unique Gabor features were extracted	162
7.13 Distance Measure D for the fitness function	165
7.14 The recognition rate versus the number Gabor feature detectors	168
7.15 Recognition rate with varying w_D	169
7.16 10 and 20 person-specific features extracted for a particular individual in the database	170
8.1 An example false face detection.	174
8.2 Block diagram.	175
8.3 Skin pixels in nRGB space.	178
8.4 Extra region for background modeling.	179
8.5 Example of <i>true</i> and <i>false</i> face detection.	180
8.6 Pre-processing.	182
8.7 Neighborhood System.	183
8.8 Frontal face Local Conditional Probability Density (LCPD) models.	184
8.9 Skin-region masks.	185
8.10 Soft threshold.	186

Figure	Page
8.11 An example of combining evidence from two experts under Dempster-Shafer Theory.	187
8.12 Coarse pose estimation.	190
9.1 Person of interest at a short distance from camera	192
9.2 Person of interest at a large distance from camera	192
9.3 Simple Background	194
9.4 Complex Background	194
9.5 Rigid, Homogeneous Object	194
9.6 Non-Rigid, Deformable, Non-Homogeneous Object	194
9.7 Static Camera	195
9.8 Mobile Camera	195
9.9 Changing Illumination, Pose Change and Blur	196
9.10 SMSPF - Step 1	201
9.11 SMSPF - Step 2	202
9.12 Structured Search	203
9.13 Sliding window of the Structured Search (Green: Estimate; Red: Sliding window).	204
9.14 Structured Search Matching Technique	206
9.15 Incorporating Chamfer Matching into Structured Search	208
9.16 SMSPF Results	210
9.17 AO (Dotted Line: Color PF; Solid Line: SMSPF)	212
9.18 DC(Dotted Line: Color PF; Solid Line: SMSPF)	213
9.19 Evaluation Measure for DataSet 1	213
9.20 Evaluation Measure for DataSet 2	214
9.21 Evaluation Measure for DataSet 3	214
10.1 (a) Typical use of the social interaction assistant, a third person perspective on the use case scenario, (b) An example of face detection being translated to vibrations on the haptic belt.	217
10.2 Main Controller implementation	229

Figure	Page
10.3 (a) 25%; (b) 75% Pulse-width modulation; (c) and (d) Vibration motor magnitudes of 25% and 75% achieved using duty cycles with 25 pulses over a 50 ms vibration period.	232
10.4 Sample Temporal Rhythm Sequences (TRS) with different magnitudes of vibration encoded on the Temporal Rhythm Units (TRU) (a) 100% Magnitude, (b) 50% Magnitude.	233
10.5 Tactor Controller implementation	233
10.6 a) Graphical User Interface on a Portable Platform. b) Temporal Rhythm Sequence (TRS) Design Interface.	234
10.7 Experiment 1 Results: Mean Localization Accuracy for each Tactor, Averaged across Subjects, with 95% Confidence Intervals	237
10.8 Mean Classification Accuracy of Duration, Averaged across Subjects and Tactors, with 95% Confidence Intervals. Durations listed in figure correspond to 200 ms (#1), 400 ms (#2), 600 ms (#3), 800 ms (#4) and 1000 ms (#5)	240
10.9 The on/off timing values of the four tactile rhythm designs, and corresponding distances, used in the experiment.	242
10.10 Overall direction recognition accuracy of each tactor location with standard deviations.	245
10.11 Overall distance recognition accuracy of each rhythm type with standard deviations.	245
10.12 Application of haptic belt as a navigational aid.	246
10.13 Arrangement of 8 Tactors around the Waist.	250
10.14 Modified Box Dance.	251
10.15 Modified Electric Slide Dance.	252
10.16 Usability Results.	256
10.17 Functionality and Performance Results.	257
10.18 Pattern Recognition Results for Dance Experiment.	258

Figure	Page
10.19 Questionnaire Results from Dance Experiment for Experienced Dance Participants (a) and Inexperienced Dance Participants (b). Responses from Q6-Q8 are excluded.	260
11.1 Group Interaction Assistance; (a) Scenario for group interaction assistance, (b) The integrated group interaction assistant.	262
11.2 Group Interaction Assistance; (a) Scenario for group interaction assistance, (b) The integrated group interaction assistant.	264
11.3 Depression and Loneliness of students plotted against stress levels in high, medium and low social skilled undergraduate students. (Please see text for the scales used for the measurement.)	270
11.4 Comparison of annual wage of the visually impaired, physically disabled and non-disabled population. (a) Compared by age group. (b) Compared by education level.	272
11.5 Comparison of personal lives of visually impaired, physically disabled and non-disabled population. (a) Number of times married. (b) Companionship status.	273
11.6 Framework for Erasable Biometrics	278
B.1 Find x and y to maximize $f(x,y)$ subject to a constraint $g(x,y) = c$	318

Chapter 1

SITUATIONAL AWARENESS IN EVERYDAY SOCIAL INTERACTIONS

People participate in social interactions every day with friends, family, co-workers and strangers. A strong set of social skills is important for a successful and productive life. For example, they help us make new friends, or make good first impressions at job interviews. Sociologists believe that social interactions are the underpinnings of our modern society, and are essential for social development and acceptance of an individual within our society. Such interpersonal interactions consists of exchanges of verbal and non-verbal communicative cues. The essence of humans, as social animals, is well exemplified in the way humans interact face-to-face with one another. Even in a brief exchange of eye gazes, humans communicate a lot of information about themselves, while assessing a lot about others around them. Though not much is spoken, plenty is always said. We still do not understand the nature of human communication and why face-to-face interactions are so significant for us.

Social interaction refers to any form of mutual communication between two individuals (dyadic interactions) or between an individual and a group (group interactions) [3]. Such communications typically involve many types of sensory and motor activities, as deemed necessary by the participants of the interaction. Social, Behavioral and Developmental Sociologists emphasize that the ability of individuals to effectively employ expressive behavior is essential for the social and interpersonal functioning of our society. Such social behaviors not only facilitate bilateral communication, but also provide a vital loop for shaping social behavior, towards developing efficient and effective social and communicative skills. Further, researchers have revealed an unconscious tendency in humans to imitate the mannerisms of their interaction partners. An increasing number of experiments have suggested that this tendency is very primeval, and that imitation plays an important role in building trust and confidence between individuals.

Unfortunately, non-verbal behavior (such as imitation) are sometime inaccessible, such as the case where the interacting participants are communicating over telephones and

not able to access the other's communicative social cues, or in the case where one or more of the interacting individuals is visually disabled, and finds it more difficult to receive non-verbal social cues.

1.1 Components of Social Interactions

From a neurological perspective, social interactions result from the complex interplay of cognition, action and perception tasks within the human brain. For example, the simple act of shaking hands involves interactions of sensory, perceptual, motor and cognitive events. Two individuals who engage in the act of shaking hands have to first make eye contact, exchange emotional desire to interact (usually happens through a complex set of face and body gestures, such as smile and increased upper body movements), determine the exact distance between themselves, move appropriately towards each other maintaining interpersonal distance that is appropriate for their cultural setting, engage in shaking hands, and finally, move apart, assuming a conversational distance, which is invariably wider than the handshake distance. Verbal exchanges may occur before, during or after the handshake itself. This example shows the need for sensory (like the visual senses of face and bodily actions, and auditory verbal exchange), perceptual (like understanding expressions, and distance between individuals), and cognitive (like recognizing the desire to interact, and engaging in verbal communication) exchange during social interactions.

Historically, social interactions have been studied in the context of human interpersonal communication dynamics under two important categories [4], namely,

- *Verbal communication*: *Explicit* communication through the use of words in the form of speech or transcript.
- *Non-verbal communication*: *Implicit* communication cues that use prosody, body kinesics, facial movements, and spatial location to communicate information that may stand alone or overlap with verbal information.

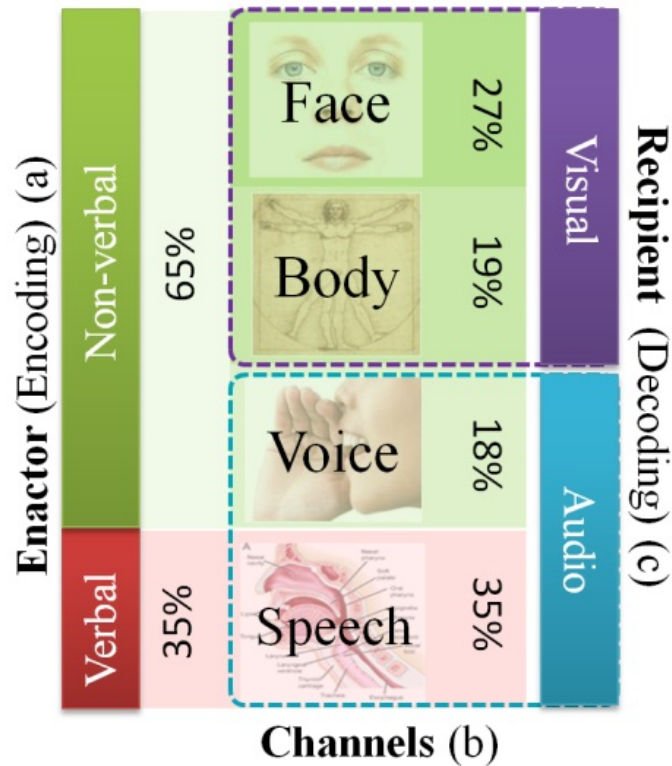


Figure 1.1: Relative importance of a) verbal vs. non-verbal cues, b) four channels of non-verbal cues, and c) visual vs. audio encoding and decoding of bilateral human interpersonal communicative cues. Based on the meta-analysis presented in [1].

1.1.1 Non-verbal communication cues

In everyday social interactions, people communicate so effortlessly through both verbal and non-verbal cues that they are not aware of the complex interplay of their voice, face and body in establishing a smooth communication channel. While the spoken language plays an important role in communication, speech accounts for only 35% of the interpersonal exchanges. Nearly 65% of all information communication happens through non-verbal cues [5]. Out of this large chunk, 48% of the communication, is through visual encoding of face and body kinesics and posture, while the rest is encoded in the prosody (intonation, pitch, pace and loudness of voice) [6]. A closer look at the various non-verbal communication modes reveals the importance of the multi-modality of social exchanges (See Figure 1.1). A component of the non-verbal cueing that is not in the figure is social touch. As will be

seen later, touch is an important aspect of social interactions that has only recently gained attention, and is now being studied extensively by behavioral psychologists.

1.1.1.1 Social Sight and Social Hearing

Unlike speech, which is mostly under the conscious control of the user, non-verbal communication channels are engaged from a subconscious level. Though people can increase their control on these channels through training, individuals are not always able to control their non-verbal cues. The unconscious revealing of one's emotional state through non-verbal channels is referred to as *leakiness* [7] and sensitive humans have learnt (possibly evolutionarily) to efficiently pick up these leaked signals during social interactions. For example, people can read very subtle body mannerisms to sense the mental state of their interaction partner. Interpretation of eye gaze is a classic example of this human ability to pick up subtle cues. Through the interaction partner's eyes, individuals can detect interest, focus, involvement, and role play, to name a few.

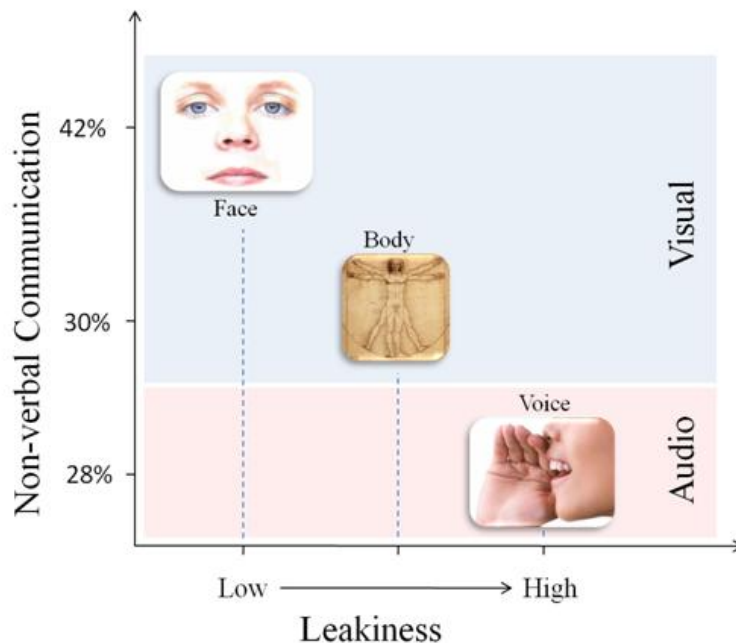


Figure 1.2: Relative communicative information plotted against its leakiness. Speech forms the verbal channel. Face, body and voice form the non-verbal communication channels.

On the leakiness scale, it has been found that the voice (not speech, but the prosody, intonation, pitch and volume) is the leakiest of all channels, implying that the emotions of people are revealed first in their voice, before any of the other channels are engaged. The voice is followed by body, face and finally the verbal channel, speech. The leakiness is plotted on the abscissa of Figure 1.2 with the ordinate showing the portion of the total information encoded in the three non-verbal communication channels. It can be seen that the face communicates the largest portion of non-verbal cues, the body communicates the next largest, and the prosody (voice) communicates the smallest, although it is the first channel to leak emotional information.

1.1.1.2 Social Touch

Apart from visual and auditory channels of social stimulation, humans rely on social touch during interpersonal interactions. For example, hand shake represents an important aspect of social communication conveying confidence, trust, dominance and other important personal and professional skills [8]. Social touch has also been studied by psychologists in the context of emotional gratification. Wetzel [9] demonstrated patron gratification effects through tipping behavior when waitresses touched their patrons. Similar studies have revealed the importance of social touch and how conscious decision making is connected deeply with the human affect system. In the recent years there has been an increased interest in the role of social touch in the enrichment of remote interactions [10] [11] in terms of an individual's social awareness and social presence.

In the next section, we discuss the concept of *Social Situational Awareness*, the importance of encoding and decoding of the social sight, hearing and touch information, and the need for individuals to be aware of their social situation for effective social communication.

1.2 Social Situational Awareness

Social Situational Awareness (SSA) is the ability of individuals to receive the visual, auditory and touch-based non-verbal cues, and respond appropriately with their voice, face

and/or body (touch and gestures). Figure 1.3 represents the concept of consuming social cues and reacting accordingly to the needs of the social situation. Social cognition bridges stimulation and reciprocation, and allows individuals to interpret and react appropriately to non-verbal cues.

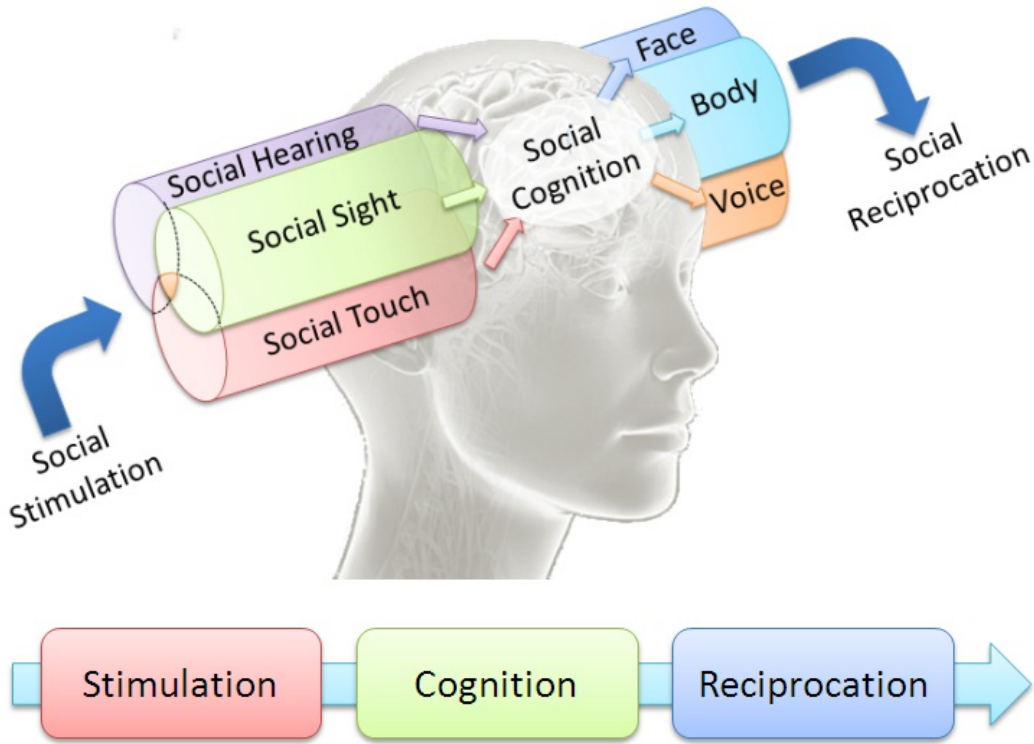


Figure 1.3: Social situational awareness in human social communications

The Transactional Communication Model [12] suggests that during any face-to-face interaction, the interpretation of the social stimulation, and the corresponding social response, are influenced by factors including culture, physical and emotional state, experience, memory, expectation, self concept, and attitude of the individuals involved in the interaction. In order to effectively process and react to the social stimulation, individuals must be able to receive and synthesize these factors. The measure of a mediating technology's ability to support social situational awareness (such as, telecommunication technology for remote interactions or social assistive technologies for the disabled population) is its ability to engage the social cognition of an individual by providing access to the above mentioned

factors and thereby evoke appropriate social reciprocation.

1.2.1 Social Situational Awareness in Everyday Social Interactions

1.2.1.1 SSA in Dyadic Interactions

Human communication theories have studied dyadic or bilateral interaction between individuals as the basis for most communication models. Theories of leadership, conflict and trust are based on dyadic interaction primitives, where the role of the various non-verbal cues is heightened, due to the face-to-face nature of dyadic interactions. Eye contact, head gestures (nod and shake), body posture (conveying dominance or submissiveness), social touch (hand shake, shoulder pat, hug, etc.), facial expressions and mannerisms (smile, surprise, inquiry, etc.), eye gestures (threatened gaze, inquisitive gaze, etc.) are some of the parameters that are studied closely in dyadic understanding of human bilateral communication [13].

✕ *Enriching SSA in mediated dyadic communication* is best done by extraction and delivery of face, body and voice-based behaviors between interacting individuals.

1.2.1.2 SSA in Group Interactions

Group dynamics refer to the interactions between members of a team assembled for a common purpose. Teams of medical professionals operating on a patient, a professional team meeting to accomplish a certain goal, or a congressional meeting to discuss regulations are all examples of groups of individuals with a shared mental model of what needs to be accomplished. Within such groups, communication behaviors play a vital role in determining the dynamics and the outcome of the meeting. Zancanaro et al. [14] and Dong et al. [15] presented one model for identifying the role played by each of the participants in a group discussion. They identified two distinct categories of roles for the individuals within the group: (1) the socio-emotion roles, and (2) the task roles. The socio-emotional roles included the protagonist, attacker, supporter and neutral, and the task roles included the orienteer, seeker, follower and giver. These roles were dependent heavily on the emotional states of the individuals participating in the group interaction. Good teams are those where

individual team members and their leaders are able to compose and coordinate their affect towards a smooth and conflict free group interaction. Effective leaders are those who can read the emotional state (i.e. affect) of each group member, make decisions about each individual's roles, and steer the group towards effective and successful decisions. Inability to assess the affective cues of team members might have significant consequences ranging from unresolved conflicts and underproductive meetings to the death of a patient.

✕ *Enriching SSA in mediated group interactions* is best done by extraction and delivery of team's interaction dynamics (as well as the group's mutual and group affect) to other participating members of the team, such as a remotely located team member or a co-located individual who is disabled.

Inadequate social awareness can lead to interactions where individuals are not engaged cognitively, and find it difficult to focus their attention on the communication. This can occur in the case of remote interactions, perceptual disabilities of some team members, and situations where medical professionals are operating simultaneously on a patient. In such cases, SSA enrichment technologies should be designed to provide a richer interaction experience for individuals involved either in dyadic or group interactions.

1.2.2 *Learning Social Awareness*

Figure 1.3 represents a simple unidirectional model of social stimulation and reciprocation. In reality, social awareness is a continuous feedback learning system where individuals are learning through prediction, reciprocation, observation, and correction of their behaviors. It is this learning mechanism that allows people to adapt their behaviors from one culture to another - here we refer to the term culture broadly, encompassing work culture, social culture in a new environment, and culture of a new team. Figure 1.4 shows the continuous feedback loop involved in social learning systems, based on the model of human cognition as proposed by Hawkins [16].

People exposed to everyday social interactions learn social skills from the three different social stimulations (social sight, social hearing, and social touch). When faced

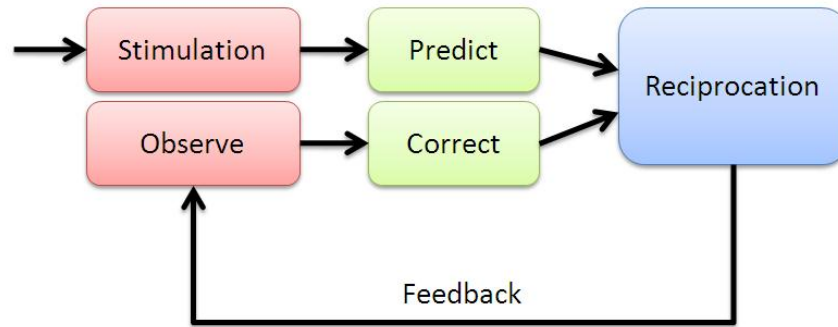


Figure 1.4: Social learning systems with continuous learning feedback loop.

with a new environment, individuals exercise their learned social skills to predict what social actions are appropriate in the new setting. Once executed, they observe and assess their counterparts, to determine if their behavior is appropriate or not for the new setting. Such learning continues until their social rule set adapts to the new environment.

Psychologists have been studying the nature of learning that happens in individuals who move from Western to Eastern cultures and vice versa. Largely, USA and Japan have been the countries of choice based on their economic equality and cultural contrasts [17]. In the West, large body movements and excitement in the voice are considered to be typical, and to a large degree are encouraged as a good social skill. Similar behavior in the East are considered to be inappropriate in professional settings and, to a large extent, indecent. An individual displaying any such inappropriate mannerisms or gestures will receive social feedback from his counterparts in the form of everyone staring at the individual, and reducing their interaction with the individual. Thus, social awareness is based on a learned set of rules about the environment within which the individual is present and this learning process requires continuous monitoring of the various social channels of stimulation. Deprivation of any one of these channels can adversely affect the ability of an individual to learn social actions and responses that are pertinent to a new social situation. Thus, enriching SSA not only offers the means for individuals to make appropriate social decisions, but also cognitively trains them towards effective social judgments.

Rest of the chapter argues that the social information impoverishment characteris-

tic of some interaction scenarios can result in social separation. Examples of such social scenarios are introduced in Chapter 2 as motivations for the development of proposed technologies that mediate social interpersonal interactions. An effort is made to identify the social separation created due to physical separation of interaction partners, and to contrast these situations with information impoverishment due to sensory/physical disabilities in co-located interaction partners. The following section highlights some of the important factors that are important to successful non-verbal communication.

1.3 Factors that are important to successful non-verbal communication

The success of non-verbal communication is influenced by various factors, some of which are dependent on the interaction partners and some of which depend on the environment where they are interacting. Psychologists have classified these factors into three categories [5]:

- (a) Factors related to the communication environment
- (b) Factors related to the physical characteristics of the communicators
- (c) Factors related to the behaviors of the communicators

Below, these three categories are each discussed in detail, providing a high level view of their influence on non-verbal communication between individuals.

1.3.1 Factors related to the Communication Environment

The communication environment (or surroundings where the interactions are taking place) makes a difference in how humans respond or react to each other [18] [19]. For example, lengthy periods of extreme heat [20] are known to increase discomfort, increase irritability, reduce work output, and produce unfavorable evaluations of others. Along with interaction partners, the environment can either have a positive or a negative influence on the emotional state of an individual. For example, wide open spaces and natural environments are known to be conducive for psychological stability [21]. Though such environmental factors are just

perceptual, they influence how humans react to each other. Some of the important environmental factors that affect interpersonal communication and non-verbal cueing are shown in the Table 1.1. This table lists eight environmental factors and provides important references from the behavioral psychology literature that discusses the influence these factors have on the non-verbal communication.

Table 1.1: Eight factors of the environment that can affect interpersonal communication.

The Communication Environment	
Familiarity of the environment	[22] [23]
Colors in the environment	[24] [25]
Other people in the environment	See next two subsections.
Architectural Designs	[26]
Objects in the environment	[27]
Sounds	[28] [29]
Lighting	[30]
Temperature	[20]

1.3.2 *Factors related to the Physical Characteristics of the communicators*

The physical appearance of a person is an important factor for non-verbal communication. People form impressions of their communication partner as soon as they engage with them. The human body communicates important sociological parameters such as status, interest, and dominance. Researchers have found both cultural and global preferences in overall body image, and any deviations from the norm affects interactions between people. For example, facial babyishness [31] has been found to affect judgment of facial attractiveness, honesty, warmth and sincerity. Deviations from the babyishness has been correlated to reductions in the judgment of these positive traits. Another example is the clothing that people wear. It has been found that first impressions are positive if the interviewer and interviewee are clothed similarly [32]. Table 1.2 shows ten important aspects of a person's physical appearance that affect interpersonal interactions. Various psychological studies have been conducted to better understand human perception of character. Although very little is known about how norms for character perception are established, the subject is being studied vigorously, especially in the context of group behaviors and personal mannerisms

within work environments [33]. Similar to Table 1.1, Table 1.2 lists the important physical characteristics that can affect communicative behaviors of interaction partners and provides references from the behavioral psychology literature.

Table 1.2: The physical characteristics of a communicator that can affect interpersonal communications.

The Physical Characteristics	
The human facial attractiveness	[31] [34] [35]
Body shape	[36] [37]
Height of a person	[38]
Self image	[39]
Body color	[40]
Body smell	[41] [42] [43]
Body hair	[44]
Clothing	[32] [45]
Personality	[46] [47]
Body decoration or artifacts	[48]

1.3.3 Factors related to the Behaviors of the Communicator

The last of the three categories of factors that affect non-verbal communication is the behavior of the communicators. The term behavior is used loosely here, as it encompasses both the static posture and the dynamic movements of the communicators. Of the three categories of factors discussed here, the behavior category is the most important. Most of the emotional information is delivered through the behavior of individuals during social interactions. Gestures, Posture, Touch, Face/Head, Eye Behaviors and Voice form the basic subdivisions in behavioral non-verbal cueing. These important aspects of non-verbal cueing are discussed below with references to various related works in the behavioral psychology that highlight the importance of these behaviors in mutual social interactions.

1.3.3.1 Gesture

Gestures are dynamic movements of the face and limbs during interpersonal communication. Together, they convey information that can be complimentary to speech, or supplementary to verbal communication. Gestures are typically classified based on their occur-

rence with speech. Accordingly, there are

- (a) *Speech-independent gestures*, or emblems (like shrug, thumbs up, victory sign etc), that are mostly visual in nature, and convey the user's response to the situation [49] [50].
- (b) *Speech-related gestures*, or illustrators (pointing to a thing, drawing a shape while describing etc) [51].
- (c) *Punctuation gestures*, that emphasize, organize, and accent important segments of a communication, like pounding the hand or raising a fist in the air.

1.3.3.2 Posture

Posture refers to the temporary limb and body positions assumed by individuals during interpersonal interactions. Posture is effective for communicating important non-verbal cues, such as leadership, dominance [52], submissiveness and social hierarchy [53]. For example, people who show a tendency toward dominance tend to extend their limbs while sitting, thereby displaying an overall larger body size. Similarly, submissiveness seems to be correlated to reducing the overall body size, by keeping the limbs together. Both gestures and postures are influenced heavily by the cultural background of the individuals, and also vary with the geographical location [54] from where they hail.

1.3.3.3 Touch

Social touch is a very important aspect of non-verbal communication in humans. Developmental biologists believe that the first set of sensory responses in a human fetus is touch [55]. From a social context this sensory channel useful for conveying interpersonal cues such as interest, intimacy, warmth, confidence, leadership and sympathy [56]. Touch is a powerful means of unconscious interaction, and people who are very good in their social skills rely upon touch a lot [57]. The sense of touch (Haptic Communication [58]) has been studied by psychologists with respect to its role as a human sensory system. More recently, haptics has been studied by technologists with respect to the role it might play in human

machine interfaces that augment or replace visual and auditory interfaces [59] with touch encoded data.

1.3.3.4 Face/Head

While the aspects of permanent facial appearance are important in the recognition of the individual, from a non-verbal communication perspective, the primary function of the face is directed towards communicating emotions and expressions. In fact face, together with the head, is the primary channel for non-verbal communication. Humans are efficient in conveying and interpreting information through subtle movements of their face and head. The facial appearance of an individual is due to (1) their genetic makeup, (2) transient moods that stimulate the facial muscles, and (3) chronically held expressions that seem to become permanent. Reliance on the face for non-verbal cues develops from a very young age. At the age of only 2 months, infants are adept in understanding facial expressions and mannerisms [60]. The human face has very fine muscular control, allowing it to display complex patterns that are widely understood among humans while being very individualistic [61]. The human visual system is able to interpret subtle differences between people's faces due to genetic makeup (i.e. person's identity through face recognition), as well as transient changes (i.e. facial expression and emotion recognition), and permanent expressions on the face (i.e. default or neutral face of individuals).

The understanding of the human facial expression space was immensely increased by the work of Ekman, Frisen [62] and Izard [63] in the late 1970s. They independently measured precise facial movement patterns, and correlated individual localized movements with facial expressions on the human face. While Izard developed these patterns on infants, the Facial Action Coding System (FACS) developed by Ekman and Frisen has become the *defacto* standard for measuring facial expressions and emotions in individuals. FACS allow researchers to encode facial movements into accurate contraction and relaxation of facial muscles. Based on these facial actions, Ekman and Frisen developed a classification of facial expressions based on six basic emotions, namely Happiness, Sadness, Anger, Disgust, Fear and Surprise. The emotions have been found to be common across cultures and age

Table 1.3: FACS communicative actions on the human face

1	Inner Brow Raiser	24	Lip Pressor
2	Outer Brow Raiser	25	Lips part
4	Brow Lowerer	26	Jaw Drop
5	Upper Lid Raiser	27	Mouth Stretch
6	Cheek Raiser	28	Lip Suck
7	Lid Tightener	29	Jaw Thrust
9	Nose Wrinkler	30	Jaw Sideways
10	Upper Lip Raiser	31	Jaw Clencher
11	Nasolabial Deepener	32	Lip Bite
12	Lip Corner Puller	33	Cheek Blow
13	Cheek Puffer	34	Cheek Puff
14	Dimpler	35	Cheek Suck
15	Lip Corner Depressor	36	Tongue Bulge
16	Lower Lip Depressor	37	Lip Wipe
17	Chin Raiser	38	Nostril Dilator
18	Lip Puckerer	39	Nostril Compressor
19	Tongue Out	41	Lid Droop
20	Lip stretcher	42	Slit
21	Neck Tightener	43	Eyes Closed
22	Lip Funneler	44	Squint
23	Lip Tightener	45	Blink
		46	Wink

levels. This fact has further motivated technologists to base human machine interaction on the detection of these emotion primitives.

The Facial Action Coding System (FACS): FACS is a systematic description of all possible facial movements in terms of muscular contractions and relaxations, as displayed by the various facial muscles, shown in Figure 1.5. There are 46 *Action Units* (AU), that form the basis of FACS. Each facial feature movement patterns is represented by 5 distinct levels of movement (A, B, C, D and E) which represent the intensity of their movement. These 46 movements are combined to represent the movements of facial features such as lips, eye brow, and chin during all communicative interactions. Table 1.3 shows the AUs that form the basis of FACS based facial coding with the appropriate number and a short description of the associated facial feature movement.

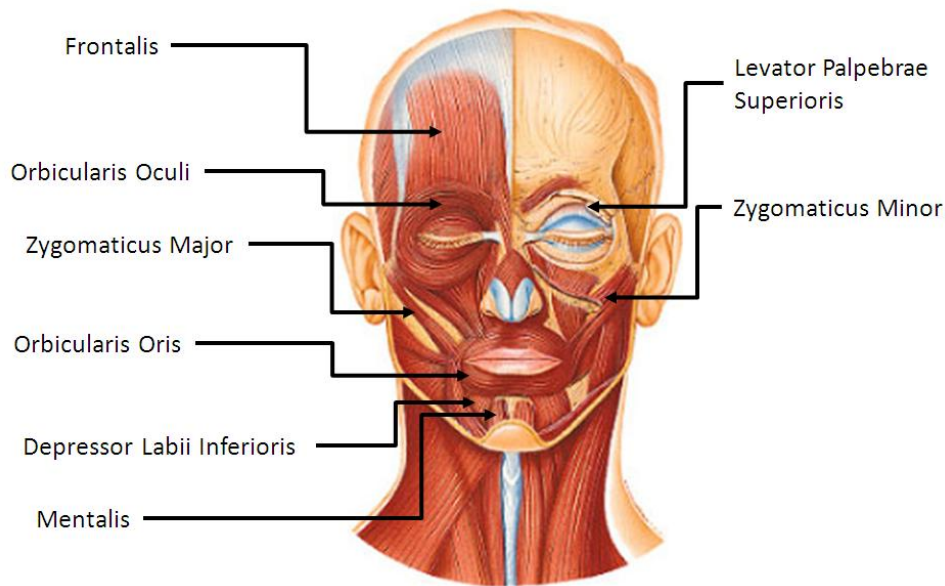


Figure 1.5: The facial muscles responsible for facial expressions and gestures.

1.3.3.5 Eye

The human eye plays an important role in non-verbal communication among interacting individuals. This involvement of human eyes is emphasized by the functions that gaze and mutual gaze play in everyday human interpersonal communication [64]. People use gaze to facilitate smooth verbal interactions that lead to information exchange [65]. The function of gaze has been classified into four important functional categories [66]. These include

1.4 Facilitating Social Interactions through the Enrichment of Social Situational Awareness

From the above discussions, it can be argued that the social interactions between individuals can be affected positively by enriching their social situational awareness. That is, by enhancing their access to social cues that could more fully engage them socially, thereby more effectively eliciting social reciprocation. The rest of this dissertation discusses the importance of social cue enrichment under various social conditions and explores one important application of human interpersonal communication enrichment in people who are

Table 1.4: The role of the eyes in human interpersonal communication.

Regulating the flow of communication	One of the most important functions of gaze is the regulation of verbal communication in bilateral and group communications. People use gaze to shift focus, bring the attention of a group of people to one thing, turn taking in group conversations [67] and eliciting responses from communication partners [68].
Monitoring feedback from the listener	Gaze provides a means for individuals to get feedback during conversations and communications. Feedback is very important while people converse. Humans study the eyes of the listener to decide whether to inject or eliminate verbal information from the conversation [69].
Reflecting cognitive activity of the communicator	Both listeners and speakers tend not to gaze at others when they are processing complex ideas or tasks. Studies have shown that people can answer better when they close their eyes and are allowed to process their thoughts [70]. Thus, cognitive processing is displayed very elegantly through eye gaze patterns.
Expressing emotions	Along with the facial muscular movements, the eyes play a vital role in the expression of emotions. In fact, in human computer interaction research, it has been found that relying on the eyes and the eyelids alone can provide more accurate affect information than relying on the entire face [71]. Verbal communication tends to move the lips and the mouth quickly and randomly. This can make image and video processing of expressions very difficult. Some of the more recent spontaneous expression recognition research is focusing exclusively on the eyes for this reason.

visually impaired. In doing so, emphasis is placed on (1) understanding the importance and priority of social signals, (2) developing sensing technologies that can extract social cues from the communicative environment, and (3) developing technologies that can deliver extracted social information to the users of the enrichment technology with a minimal sensory and cognitive load.

1.5 Organization of the Dissertation

The dissertation is organized as shown in the Table 1.5

Table 1.5: Organization of the dissertation.

Chapter 2	Discusses the need for enriching social situational awareness in everyday personal and professional lives of individuals.
Chapter 3	Highlights the importance of enriching social situational awareness for individuals who are blind and visually impaired and lays foundation for the bulk of the work presented in this dissertation.
Chapter 4, 5, 7 & 8	Discusses various technologies that can enable users who are blind and visually impaired to access social signals that are important for having a rewarding social interaction with sighted counterparts. The details of these chapters will become clear once the reader has been introduced to the various social situations that require attention, as detailed in Chapter 3.
Chapter 6 & 10	Discusses various technologies that can enable any processed social signals to be delivered to people who are blind and visually impaired, without overloading any of their senses, like hearing or touch.
Chapter 11	In wake of introducing technologies that interject into the personal lives of individuals (lives of not only those who are disabled but also their interacting partners) this chapter highlights the need for researchers to consider the impact of social mediation technologies on the society. While the discussions do not impart strict policies, this chapter initiates a conversation towards adopting important technology policies in the emerging assistive technology domains.

Chapter 2

NEED FOR ENRICHING SOCIAL SITUATIONAL AWARENESS

Social situational awareness (as described in Figure 1.3) can have a tremendous impact on the quality of one's personal and professional life as well as one's ability to relate to their social surroundings. This chapter introduces three important example situations where there is an increased need to enrich social situational awareness.

1. where a person with a disability is interacting in a social situation with sighted counterparts.
2. where two or more people are interacting from remote locations.
3. where medical teams are interacting to perform an operation on a patient.

2.1 Disability Induced Social Signal Attrition

Due to the fact that a large portion of human-human interpersonal communication happens through complex non-verbal cueing, individuals who are disabled face myriad levels of difficulty when it comes to interpreting and responding to everyday social interactions. The difficulty varies based on the kind of disability and the intensity of the disability one faces. Non-verbal cues are mostly interpretative and not instructive as verbal cues (such as speech) are. In a bilateral interpersonal interaction, while speech encodes all the information, non-verbal cues facilitate an elegant means for delivery, interpretation and exchange of the verbal information. People with sensory, perceptive, motor and cognitive disabilities may not be able to receive or process these non-verbal cues effectively. Though most individuals learn to make accommodations for the lack of a primary information channel, and lead a healthy personal and professional life, the path towards learning effective accommodations could be positively effected if social signals could be enriched for the benefit of these individuals. We focus on the topic of building technologies that can mediate interpersonal interactions for people who are disabled and we specifically focus on the issues

emanating from the lack of *sensory visual channel*, like in the case of people who are blind or visually impaired.

2.1.1 Visual Impairment - a hinderance to smooth social interactions

As seen in Figure 1.1, most non-verbal cues are perceived visually. While some of these cues are delivered along with speech, nonverbal cues such as posture or gestures is inaccessible to someone with a significant visual impairment or blindness. This deprives these individuals of vital communicative cues that normally enrich the experience of social interactions, and puts them at a disadvantage in daily social encounters. For example, during a group conversation it is common for a question to be directed to an individual without using his or her name - instead, the gaze of the questioner indicates to whom the question is directed. In such situations, people who are blind find it difficult to know when to respond because they cannot determine the direction of the questioner's gaze. Consequently, individuals who are blind might be slow to respond or they might talk out of turn, possibly interrupting the conversation. This can lead to isolation and reduced sense of engagement with the ongoing interactions.

Compounding these problems, sighted individuals are often unaware of the role their own non-verbal cues in their social interactions, and they often fail to make appropriate adjustments when communicating with people who are blind. When people who are blind find themselves in such a social interaction, they might be reluctant to ask the sighted person to make adjustments to their disability, because they do not want to burden friends and family. The combination of all these factors can lead people who are blind to become somewhat socially isolated [72]. Ironically, while people who are blind and visually impaired face difficulties in social interactions, research in rehabilitation training for these populations suggests that these individuals need to increase their social interaction to improve their acceptance in mainstream careers.

Recently, Jindal-Snape [73] [74] [75] carried out extensive research in understanding social skill development in blind and visually impaired children. She has studied in-

dividual children who are blind from India, where the socio-economic conditions do not provide for trained professionals to work with children who have disabilities. Her seminal work in understanding the social needs of children who are blind have revealed two important aspects of visual impairment that prevent social interactions. These include:

- (i) The inability to learn social skills due to the lack of vision.
- (ii) The lack of reinforcement feedback on one's mannerisms.

2.1.1.1 Inability to learn social skills due to the lack of vision:

Jindal-Snape observed that significant others in the environment often fail to give verbal replacement for their bodily actions, and even when they do, it is not meaningful or understandable to an individual who is visually impaired - for example, nodding one's head in reply to a question or gesturing. Lack of meaningful verbalizations could make it difficult for visually impaired persons to comprehend a conversation [74] [75] and, at times, may stop conversing. Similar studies carried out by Celeste [76] indicated that social intervention by parents and teachers are very important in the formative years of a child with visual impairment. Developing on the work by [77], which emphasizes that short-term feedbacks are never effective, Celeste insists that professionals must identify strategies related to social skills that work, provide consistent support and follow children longitudinally to ensure effective development of social skill set.

People who are sighted do not necessarily have the training to interact with individuals who are blind or visually impaired. Thus, unconsciously they tend to neglect people who are blind. For example, sighted people use eye contact as a primary means of keeping the attention of people they communicate with. While conversing with a person who is blind or visually impaired, sighted individuals expect eye contact. The lack of such a feedback distracts the sighted individuals to assume wandering attention or disinterest from the visually impaired individual. Research indicates that blind individuals with the ability to accommodate the expectations of their sighted counterparts have great potential for personal and professional growth.

2.1.1.2 Lack of visual reinforcement feedback on one's mannerisms:

When individuals display behaviors in any social setting, they receive feedback through their peer's reactions. These reactions could either reinforce or dissuade those behaviors depending upon whether they are deemed appropriate or not, respectively. Due to their inability to receive visual feedback, people who are blind or visually impaired do not receive this feedback from their social counterparts. People who are visually impaired at a very young age are at a particular disadvantage in learning appropriate social actions and mannerisms. Development of asocial stereotypic body mannerisms are one such case where positive reinforcement through visual stimulation is necessary to cull certain developmental behaviors (such as body rocking) that would have otherwise weaned off gradually as the child gets into adulthood.

Most people who are blind or visually impaired eventually learn to partially compensate for the lack of visual cues by using other cues, such as audio. It maybe possible that the path towards learning such accommodations could be positively effected through the use of technology that mediates interpersonal interactions. Specifically, children with visual disabilities find it very difficult to learn social skills while growing amongst sighted peers, thus avoiding social isolation and psychological problems [73].

Social disconnect due to visual disability has also been observed at the college level, where students start to learn professional skills, and independent living skills. Any assistive technology that can enrich interpersonal social interactions could prove beneficial to the visually disabled or blind people during this learning process. Technology specialist Shinohara [78] [79], observed the everyday activities of a college student named Sara who was blind [80]. Shinohara categorized Sara's daily needs into functional categories and identified 5 important aspects of Sara's life where she could benefit from assistance. These include (in order of importance)

- increased socialization.

- increased independence in doing things.
- increased control over things she does.
- feedback from objects around her.
- increased efficiency in her activities.

As seen from the list, socialization was a very important aspect of this college student's needs. Shinohara concluded that technology that support socialization for people with visual impairment is absolutely necessary.

2.2 Social Signal Attrition during Remote Interpersonal Interactions

The globalization of economies has required that people communicate across geographical distances efficiently and effectively in real (or near-real) time. This has increased international and inter-cultural interactions where people with different cultures are working together to accomplish common tasks. Intercultural interactions cause socio-emotional stress due to differences in work ethics, communication protocols and the role of hierarchy in management. Such stress is typically managed well through dyadic face-to-face interactions, which is known to reduce cultural stress and communicative misunderstandings [81]. Unfortunately, face-to-face dyadic interactions are not always possible, and people might need to rely on telecommunication technologies to bridge geographical separation. Existing technology solutions that are closest to simulating face-to-face interactions (such as telepresence environments) are typically limited to more scheduled, highly structured, and formal interactions. Furthermore, all current telecommunication technologies that support virtual collaborations suffer from emotional impoverishment, due to the limited social situational awareness among participating members of the interaction. In this section, we highlight some of the important problems faced by remotely distributed teams when the interpersonal communication is restricted to virtual telecommunications.

Kock et al. [82] expands on various studies in the area of professional communication to elaborate on the various problems faced in enriching remote interactions between

geographically isolated individuals. In doing so, they identify several important challenges for *e-collaborations*, including

- *Theoretical Challenges*: Lack of theoretical understanding of media-based human-human communications.
- *Human/User Challenges*: Lack of understanding of the human dynamics that make interpersonal interactions so important to humans.
- *Technical Challenges*: Lack of technology to provide seamless social presence across geographical boundaries.
- *Conceptual Challenges*: Shear complexity of the problem given the above three challenges.

Discussions in Chapter 3 will discuss the similarities between the above four challenges encountered when mediating remote person-to-person interactions, and the challenges encountered when developing mediating assistive technologies for people who are blind or visually impaired.

Most remote collaboration technologies rely on media richness to compensate for the lack of social presence. Early behavioral psychology studies supported this approach with Media Richness Theory [81]. To provide media richness despite communication bandwidth limitations, researchers employed avatar-based communication. However, this approach was not well-received and subsequent studies have shown the need for media naturalness in telecommunication, more so than just media richness [83]. The naturalness here refers to the sensitivity of the human communicative system to the subtle movements and gestures shown by the face, body and head during social interactions. In the context of naturalness, it is important to cite Robert et al. [84], who addressed the question of how much naturalness is necessary to involve the participants effectively in the interaction. They addressed notions that *more social presence in media is always better* and through experiments they theorize that, “*the use of media high in social presence induces increased*

motivation but decreased ability to process information, while the use of media low in social presence induces decreased motivation but increased information processing ability.”

Their work in the potential negative consequence of high social presence further increases the need to understand how social interactions across large geographical boundaries can be bridged. In this dissertation, we argue that the social engagement of individuals can be enriched through the use of novel multimodal social situational awareness technologies that prevent overloading of any one sensory channel of the participants and encourages distribution of the social signal delivery across various under utilized human sensory channels.

An industry survey [85] of 1592 individuals who collaborated remotely, carried out by RW3 CultureWizard - a company focused on improving international collaborations - reported difficulties that are representative of (1) lack of commitment in remote interactions and (2) inability of virtual teams to correspond and communicate as effectively as face-to-face teams. “Respondents found virtual teams more challenging than face-to-face teams in managing conflict (73%), making decisions (69%), and expressing opinions (64%). The top five challenges faced during virtual team meetings were insufficient time to build relationships (90%), speed of decision making (80%), different leadership styles (77%), method of decision making (76%), and colleagues who do not participate (75%).” These results can suggest a need for Social Situational Awareness in group settings. In the classical model for group dynamics, Bruce Tuckman [86], defines four stages in the formation of an efficient group. *Forming*, *Storming*, *Norming* and *Performing* describe the typical process that the groups go through before delivering at their best. The stages of *Storming* and *Norming* are deeply connected to the individual group member’s abilities to communicate, coordinate and empathize with their fellow group members. The socio-emotional interactions between the group members dictate how quickly (or slowly) a group will progress from the *Formative* first stage to *Performing* fourth stage [87].

Further, when the participants were asked about the personal challenges faced during virtual team meetings, they reported inability to read non-verbal cues (94%), absence of collegiality (85%), difficulty establishing rapport and trust (81%), difficulty seeing the

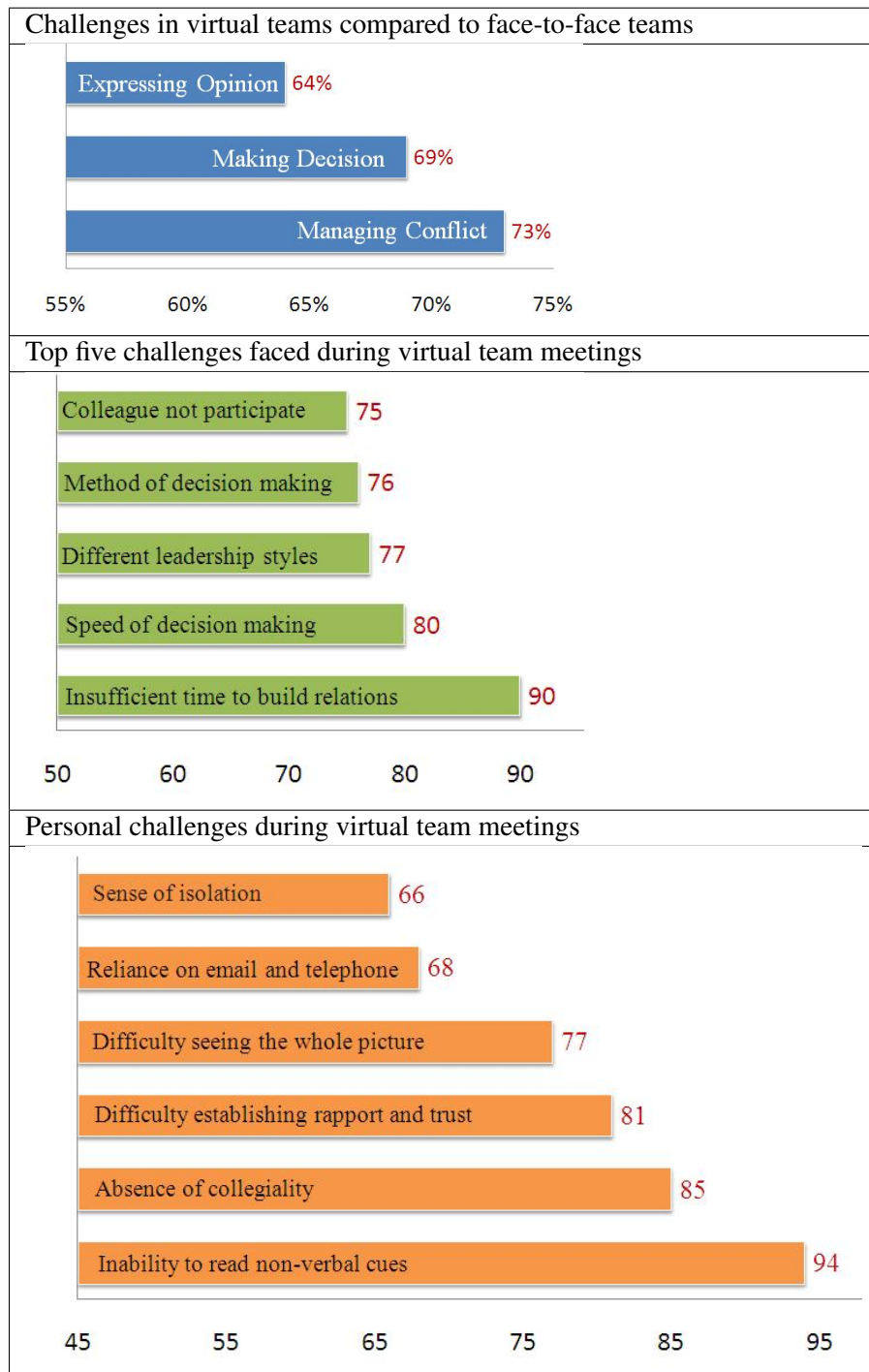
whole picture (77%), reliance on email and telephone (68%), and a sense of isolation (66%).” Delivering non-verbal cues, establishing trust and rapport, and easing isolation are all factors that are important for increasing one’s social connection to their interaction partners, be it remote or face-to-face. This result is in accordance with psychology studies carried out on e-collaborations. The review of Media Richness and Social Presence theories by Kock [88] highlights problems that are very similar to those described by the CultureWizard studies.

As seen from the discussions above, remote social interactions rely on the social presence and situational awareness of the interacting participants. Enriching any of the social cues (sight, hearing or touch) through remote means could potentially improve communication dynamics. As will be discussed in Chapter 6, various attempts have been made to mediate complex interpersonal interaction signals across physical separation. In the following section, we discuss an interesting social interpersonal communication issue that arises within a professional environments, specifically critical care medical teams, and show how social situational awareness is of utmost importance.

2.3 Social Signal Attrition in Medical Teams

Modern day critical care facilities require the simultaneous efforts of multi-disciplinary medical professionals (such as doctors, surgeons, nurses and anesthesiologists) to operate on a single patient. This imposes a mandatory requirement on the professionals to work as a team. Unlike many other professional teams who choose their members after careful deliberation, medical teams are assembled dynamically, based on whichever medical professionals are available on the hospital floor at the time of emergency. Further, these teams might last for a very short duration of time (the duration dictated by the emergency) dissolving after the need has been met with different teams forming for subsequent emergencies. Studies show that teams that establish well articulated communication between members perform well under the stressful environment. Unfortunately, this is not true of all medical teams and one individual’s stress may very well propagate through the team, breaking down mutual communication and support, leading to the patient’s death.

Table 2.1: Survey on the challenges of remote interaction.



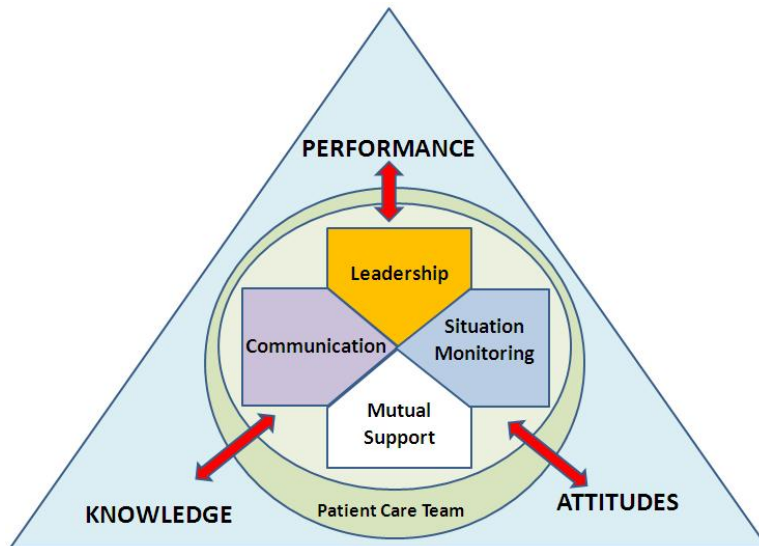


Figure 2.1: TeamSTEPPS: Team Strategies and Tools to Enhance Performance and Patient Safety

The importance of studying a group of physicians entering a medical situation as a single operating unit began to appear in the focus of behavioral scientists with the publication of the Institute of Medicine (IoM) report titled *To Err is Human: Building a Safer Health System* in Dec 1999. One of the four core messages from the report identified that patient life is lost not because of the failure of an individual, but due to the failure of the team. Since this report, Agency for Healthcare Research and Quality (AHRQ) and Department of Defense (DoD) have focused on team failure from an Evidence-Based Medicine perspective and released Team Strategies and Tools to Enhance Performance and Patient Safety (TeamSTEPPS) [89] as the standard for team training in health care. The core of TeamSTEPPS, designed to be a training tool that promotes teamwork among the multi-disciplinary members of a dynamically-formed medical team, focuses on the need to have four important team attributes among the members, namely,

- (a) Leadership
- (b) Mutual support
- (c) Communication

(d) Situation monitoring (or a shared mental model)

As seen from Figure 2.1, TeamSTEPPS focuses on the individual physician and the team's ability to work together as a system. Leadership, Communication, Situation Monitoring and Mutual Support were all derived from earlier DoD and AHRQ lead studies in medical team management and are based on the underlying principles of: Team Leadership [90], Mutual Performance Monitoring [91], Backup Behavior [92], Adaptability [93], Team/Collective Orientation [94], Shared Mental Model [89] [93], Mutual Trust [95] and Closed loop Communication [92] (For a detailed analysis of each of these principles, please see King et al. [96].) Most of these principles are in turn derivatives of the social skill set of the individuals who make up the medical team that is responsible for the patient safety. It has been shown that, in cases of medical errors leading to loss of life, communication breakdown between one or more team members resulted in an avalanche of problems, eventually resulting in death.

2.3.1 Importance factors affecting team performance

2.3.1.1 Group Dynamics

As discussed briefly in Section 1.2.1.2 of Chapter 1 Group Dynamics focuses on studying the various components of group interactions, including inter-agent communication [97], productivity of a given group [98], level of understanding of each other's potentials and limitations, job satisfaction and combined creativity of a team [99] to name a few. In recent years, interest in understanding group dynamics in work environments has tremendously increased in interdisciplinary teams involving computer scientists and socio-behavioral psychologists in the area of Computer Supported Collaborative Work (CSCW) [100]. In the context of medical teams, group dynamics focuses on the ability of the physicians and specialists to intercommunicate their needs. During emergency situations, group dynamics facilitates the emergence of a Shared Mental Model [101], which enables all the professionals to relate to each other in terms of what needs to be done towards resuscitating the patient.

2.3.1.2 Leadership

Theories of leadership have proposed evidence-based models for explaining qualities exemplified in successful leaders. From bureaucratic leaders to political leaders, the models used to explain the qualities of leaders vary dramatically. There is no single accepted definition of what a leader should represent, as the problem of identifying a leader is highly contextual in nature. Recently, the functional model of leadership has been developed to describe team leaders as having *self regulation* which translates to learning, performance and adaptability. These models allow the study of dynamic teams that are formed for very short durations (such as medical response teams) and allow monitoring of each individual and their contribution to the team activity [102]. Kozlowski et al. have described a dynamic multi-goal model for team leadership as shown in Figure 2.2. Based on this model, they describe effective leaders as those who can not only assess simultaneously their own goals while keeping track of team goals in a dynamically evolving situation.

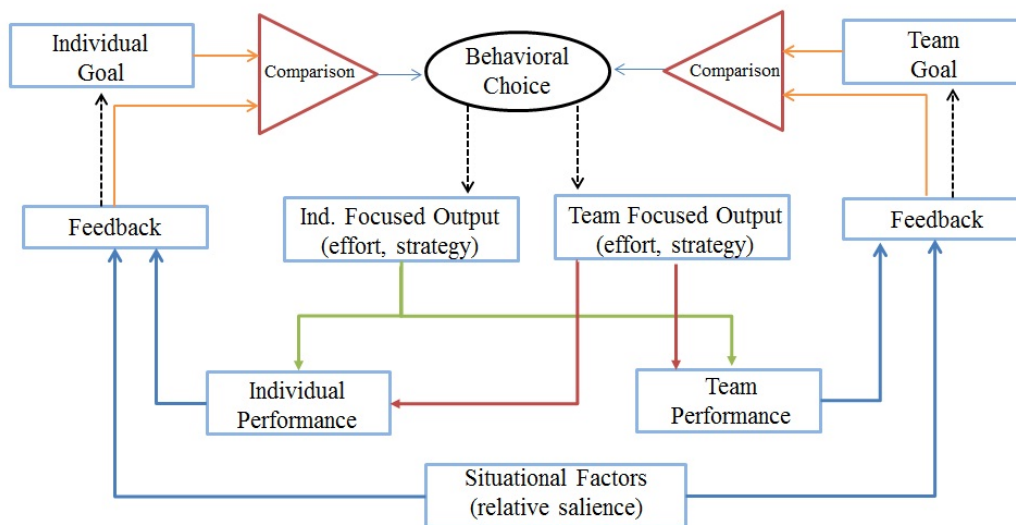


Figure 2.2: Multi goal model of self regulation for effective team leadership.

While Figure 2.2 shows the behavioral choice of the leader to be a vital component of self regulation and team regulation, very little study has been focused on the effect of leader's socio-emotional state (assessed through their social situational awareness) on team

dynamics. Recently Sy et al. [103] have demonstrated how important it is for the leader to control and regulate his/her affect cues within dynamically formed teams. The mood of the leader propagates through the team, and can have net positive or negative effect on the team outlook and performance. The dynamic nature of team formation is further complicated in medical teams as the responsibility shifts very quickly from one specialist to another, as they operate on the patient [104].

Behavioral psychologists have been studying the impact of socio-emotional states (especially stress induced socio-emotional states) on the performance of professionals and conclude that the direct artifacts of stress include deprecated decision making, failure in leadership and breakdown of mutual support. Further research shows that enriching the social situational awareness of the professionals will reduce stress and potentially improve team performance. Further, assessing the socio-emotional and communication skills of the professionals within the critical care unit will provide an unfettered advantage towards determining metrics of team performance under stress. To this end, three parameters have been identified as important towards advancing patient safety in critical care environments:

1. Automated monitoring of group dynamics to determine communication breakdowns.
2. Automatic evaluation of the social affinity between team members.
3. Leadership evaluation and nomination through long term monitoring of individuals.

2.3.2 The emerging science of medical team social assessment

In light of the important issues related to group interactions and leadership, Krishna et al. [105] [106] have suggested a study of three important challenges of medical teams through critical care simulation as discussed below. They propose to study interdisciplinary medical residents as they respond to simulations of critical care scenarios. Teams before and after TeamSTEPPS training will be assessed based on their leadership, mutual support, communication and situation monitoring. Below, the simulation facility and the proposed

research directions are discussed in detail, highlighting the importance of social situational awareness among the medical team members.

2.3.2.1 The Mayo Clinic Multidisciplinary Simulation Center:

Mayo Clinic has invested in the development of a state-of-the-art Multidisciplinary Simulation Center providing many advantages in learning team performance and hospital emergency code training. The center was designed to reflect the complexity of the live clinical environment. The physical, electronic, microcultural and macrocultural environments have been replicated in order to create a practice field for clinical situations. The physical environment includes artifacts such as the bed, sink, cabinets, flooring, doors and lights. The electronic environment replicated with monitors, intravenous pumps, electronic medical records, pharmaceutical dispensing system, decision support software and computers. The micro cultural environment is provided by the inter professional team, and the macro cultural environment is provided by embedding the simulation center within the hospital. The patient clinical situations are replicated by combinations of high fidelity simulators, virtual reality simulators, standardized patient actors, and low fidelity simulators. Experts have developed standardized simulation scenarios, curriculum and debriefing specific to learner or team performance levels.

Within this environment, a high resolution audio and video capture system has been unobtrusively placed for performance monitoring and archiving. Data is saved to an internet-based learning management system with individual and team portfolios permitting immediate local or asynchronous remote review and feedback. Teams practice rare or life threatening event management in an error-forgiving environment gaining experiences which have been shown to improve patient safety [107]. Some of the highlights of simulated team training include,

- A hospital based center, with clinical microsystems reflecting the live environment
- High fidelity simulators, virtual reality trainers, simulated electronic medical record

- Wireless biometric monitoring, an audio-video architecture for simulation suites, in situ, and in vivo performance archive
- An internet-based learning management system, with individual and interprofessional team portfolios
- Expert developed, standardized curriculum and debriefing
- Curriculum customizable to the learner or team performance level
- An error-forgiving clinical experience, enhancing patient safety
- Deliberate practice with supervised instruction in life-threatening event management
- Experience with uncommon scenarios
- Leadership training and debriefing using “Crisis Resource Management” principles

2.3.2.2 Challenges in Social Situational Assessment and Training

Challenge 1: Automated monitoring of group dynamics to determine communication breakdowns: Current team performance analysis systems are mostly based on retrospective video stream analysis collected during simulations of hospital emergency codes. The analyses are mostly based on expert opinions of what happened during the critical incidents of the simulation [108] [109]. Unfortunately, expert’s time is very valuable and post-simulation analyses may not get sufficient attention, due to heavy hospital loads. For a long time researchers have questioned how communication between medical team members vary over the period of the emergency code execution [110]. However, very little is understood about the basics of the communication patterns during emergency situations, mostly due to the lack of automated annotation systems that do not require expensive specialist time. Automated team performance analysis systems that focus on detecting specific instances of communication breakdown occurring during emergency code simulation could greatly enhance our understanding of team work and how it can be enhanced through training.

Challenge 2: Automatic evaluation of the social affinity between team members: Sociograms (social affinity maps) have been used historically to determine the interpersonal match between members of a team or an organization. Sociograms are obtained through the process of sociometry [111], which quantitatively measures the relationships of individuals who exist within a social space. As mentioned earlier, in medical teams, the social space happens to be the emergency room where the team assembles with very little or no time to determine who are the members of the team. Sociometry is achieved through a set of evaluations that can assess the social interactions between individuals. The measurements could be done within the environment where the individuals interact (the medical team) or outside (casual interactions). Technologies developed to assess sociometric affinity between professionals could in turn provide quantitative evaluations of the social interactions between individuals. Sociograms developed at a hospital level could offer effective tools for quick team formations. Teams formed out of specialists, technicians and nurses who are closer to one another on the sociogram could offer a team with relatively less emotional stress. Socially closer individuals will also exhibit better communication, thereby increasing team performance.

Challenge 3: Leadership evaluation and nomination through long term monitoring of individuals: Theories of leadership have proposed evidence-based models for explaining qualities exemplified in successful leaders. Recently, the functional model of leadership [112] has been developed to describe team leaders as having self regulation which translates to learning, performance and adaptability. These models allow the study of dynamic teams that are formed in very short durations (like medical response teams) and allow monitoring of each individual and their contribution to the group activity. Kozlowski et al. [112] also describe a dynamic multi-goal model for team leadership which models effective leaders as those who can not only assess their own goals but also keep track of team goals, while approaching a dynamically evolving situation. Technologies developed towards understanding and modeling human interactions and communications can provide the tools needed to measure leadership qualities through long-term monitoring.

In summary, social situational awareness is an important aspect for medical pro-

professionals entering into dynamic team settings requiring effective communication between team members. Enriching the awareness of one or more of its members will directly reflect on their performance as a team and hence the patient safety.

2.4 The Research Focus of this Dissertation

From the above sections, it is evident that enriching social situational awareness is an essential component of enriching interpersonal interactions for both the personal and the professional lives of individuals. Further, as can be seen from the three distinct discussions on the need for social interpersonal communication, there is no single underlying theory of social interactions and social situational awareness that can be immediately leveraged towards developing models of enrichment. Years of research into human-human behavioral dynamics have resulted in theories that are only grounded in coarse human communication experiments. It has proven very difficult to finely define the nuances of human interactions, especially in dynamic contexts. Modeling the complex nature of human interpersonal communications under all contexts of interaction is a grand challenge. Kock et al. [113] offer this opinion through an evolutionary model for human electronic communications where they argue that humans evolve continuously within their cultural contexts and no one theory can describe the inner workings completely. Kock in his seminal discussion “The Ape that Used Email: Understanding E-communication Behavior through Evolution Theory” [88] argues that a) Media Naturalness, b) Innate Schema Similarity and, c) Learned Schema Variety, are the foundations for human communication, where, media naturalness refers to the ability of humans to express subtleties in a natural way through a medium of remote communication, innate schema similarity is the similarity in communication patterns shared by members of a common culture, while learned schema variety refers to the individual differences that makes each person learn as they coexist within a culture. It is this similarity versus variety that makes each person both a member of the culture, yet an individual on their own.

Acknowledging this complexity, the research highlighted in this dissertation assumes an *Evidence-based modeling approach to enriching human-human interpersonal in-*

teractions through multimedia mediation. As the name suggests, Evidence-based modeling is based on observations of typical attributes of the problems at hand, and proposes to address those problems in a specified context, with little or no generalization beyond the boundaries of that context. Evidence-based models are common in medical practices where the outcomes of prescribed health care are not considered to be objectively measured due to individual patient factors, and can never be proven rigorously through scientific process. A popular remark [114] on Evidence-based Medicine (EBM) or Evidence-based Practice (EBP) reads,

EBM/EBP recognizes that many aspects of health care depend on individual factors such as quality-of-life and value-of-life judgments, which are only partially subject to scientific methods.

Though it has pitfalls, Evidence-based Medicine has been successful in saving lives [114] and increasing the perceived quality of life.

Based on the above observations, the research described in the rest of the dissertation follows an evidence-based method aimed at enriching interpersonal interactions among individuals. Further, this dissertation focuses on the first area introduced in this chapter, mediating interpersonal interactions for people who are blind and visually impaired.

2.4.1 The Handshake Example: An Example of Evidence-based Understanding of Social Situational Awareness

In Section 1.1 of Chapter 2, the handshake scenario was introduced as an example of the complexity involved with even simple social interactions. Here we reintroduce the handshake and describe how it presents challenges to people who are blind and visually impaired.

What seems to be a simple act of shaking hands between two individuals is actually a rather complex interplay of cognitive sensorimotor events. Two individuals who engage in shaking hands first make eye contact, signal a desire to interact (usually through face and body gestures, such as smile and increased upper body movements), determine

the necessary distance between themselves, move appropriately towards each other while respecting the interpersonal distance factors relevant to their particular cultural setting, engage in shaking hands, and finally, move apart, assuming a conversational distance which is invariably greater than the handshake distance. Verbal exchanges might or might not occur before, during or after the hand shake itself. This example shows the need for sensory (visual senses of face and bodily actions, as well as auditory verbal exchange), perceptual (understanding expressions, and social distance between individuals), and cognitive (recognizing the desire to interact, and engaging in verbal communication) activity during everyday social interactions.

People who are blind or visually impaired face numerous challenges when it comes to interactions such as handshake. They are not able to process the visual cues of where someone is standing with respect to themselves (especially in a group setting), they cannot determine if anyone has made eye contact with them (indicating a desire to interact), and they may not be able to determine exactly how far their interaction partners are located, and in what direction. As a result, they typically initiate a handshake by standing in place and extending their arm in a handshake posture in the direction where they believe their interaction partner to be, hoping to draw the attention of their sighted counterparts. In dyadic interactions, this strategy is likely to elicit a handshake. However, when there is a group of sighted individuals who are all interacting among themselves, this strategy might cause momentary confusion, as members of the group attempt to resolve uncertainty about who should respond.

Figure 2.3 represents an evidence-based model of delivering social situational awareness to an individual who is blind such that he/she can carry out a handshake social interaction amidst a group of sighted individuals. Note that the Proxemics Model presented in the figure refers to the interpersonal spaces that people occupy on a day-to-day basis and it is heavily influenced by the culture in which one resides [115]. We plan to develop and evaluate such evidence-based models as guides for the development of assistive technologies that can communicate important non-verbal cues. Chapter 3 discusses in detail how

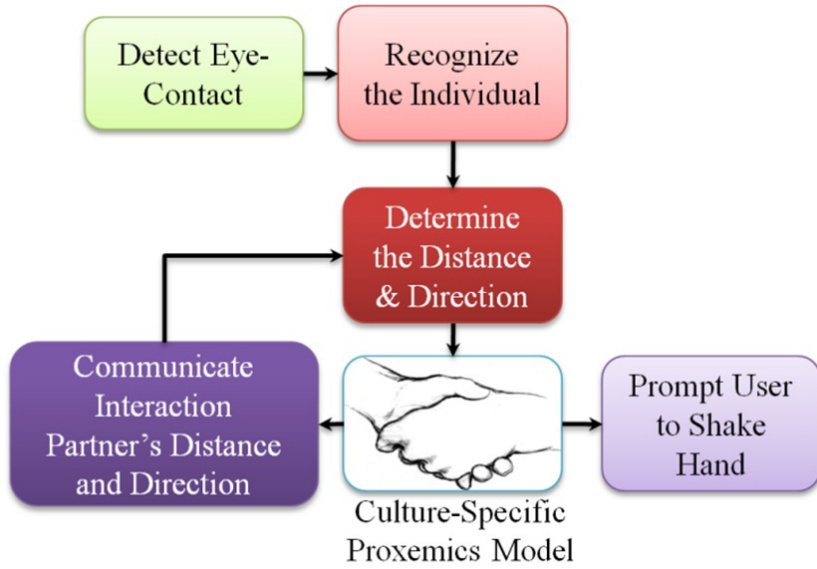


Figure 2.3: Evidence-based Model for Social Situation Awareness to promote handshake non-verbal cueing in the visually impaired and blind population.

evidence-based models of social situations for people who are blind and visually impaired were constructed for use in guiding the development of mediating technologies.

Chapter 3

ASSISTIVE MEDIATION TECHNOLOGY FOR INDIVIDUALS WHO ARE VISUALLY IMPAIRED

The overarching problem of disability-induced social signal deprivation was described in detail in Section 2.1. Further, the adverse effects of the lack of social situational awareness for people who are visually impaired were highlighted as,

- the inability to learn social skills due to the lack of visual feedback
- the lack of reinforcement visual feedback on one's mannerisms

This chapter discusses research in the related areas of a) computer vision based non-verbal cue sensing, b) social signal processing, c) human-computer interfaces, d) multimedia technology for computer human interactions and e) assistive technology design and development. The results of this research inform evidence-based models for social signal enrichment. Also presented are the results of a survey aimed at better understanding the problems faced by people who are visually impaired as they engage in social interactions with their sighted counterparts, which provides a basis for developing assistive technology for mediating social interactions.

3.1 Automatic Detection of Non-verbal Cues and Observations

Affective Computing research has quantitatively studied both verbal and non-verbal cues displayed by the humans during social communication. Signal streams from various sensors, including visual sensors (e.g. cameras), audio sensors (e.g. microphones) and various physiological sensors (such as EEG, EMG, and galvanic skin resistance sensors) have been used to sense human emotional states. A good review of research work in Affective Computing can be found in [116]. This research has enabled a better understanding of human physiological signals, with respect to emotional states, and the results have been used to facilitate human-computer interaction (HCI). In theory, a system that can remotely detect

non-verbal social cues could also be used as an assistive device to provide social feedback to people with disabilities. The emphasis here would not be so much on interpreting these cues as on presenting raw social cue data to the user, and allowing the user to interpret them. However, very little research has been done towards finding intuitive methods for presenting raw social cue data to humans. [117] developed a haptic chair for presenting facial expression information. It was equipped with vibrotactile actuators on the back of the chair that represented specific facial features. Experiments conducted by the researchers showed that people were able to distinguish between six basic emotions. However, this solution had the obvious limitation that the user needed to be sitting in the chair to use the system.

Observation 1: Assistive technology designed to facilitate social interaction should be portable and wearable, so that the users can employ them in various social circumstances without imposing significant restrictions on their activities.

People with disabilities (like blindness or Autism) are not always able to perceive or interpret implicit social feedback to guide them in improving their communication competence. However, they might be able to use explicit feedback provided by a technological device. Rana and Picard [118] developed a device called Self Cam, which provides explicit feedback to people with Autism Spectrum Disorder (ASD). The system employs a wearable, self-directed camera that is supported on the user's shoulder to capture their own facial expressions. The system attempts to categorize the facial expressions of the user during social interactions to evaluate the social interaction performance of the ASD user. Unfortunately, the technology itself plays a role in the social interactions in which they are used, becoming more of a social distraction for both participants than an aid.

Observation 2: Assistive technology designed to facilitate social interactions should be implemented in such a way that it does not itself become a social distraction.

Vinciarelli et al. [119] describe the use of technologies for understanding social interactions within groups, specifically targeting professional environments where individuals make decisions as a group. They analyze body mannerisms and prosody to extract nonverbal cues that they use to do group dynamics analysis. They rely on simple sensors

in the form of wearable tags [120] (which detect face-to-face interaction events) along with prosody analysis to determine turn taking, emotional state of the speaker, and distance to an individual. Pentland describes these signals captured during group interactions as [121] honest signals. Some of his recent work [122] in the area of social monitoring attempts to capture these signals and provide feedback to individuals about their social role within a group. The use of social feedback is illustrated elegantly in their work but their approach required all members of the group to carry sensors all through the study. This proved to be a viable and productive approach for studying group dynamics, but is not viable as a strategy for the implementation of an assistive technology, as it is not realistic to assume that everyone who interacts with a person who has a disability will wear sensors.

Observation 3: Assistive technology designed to facilitate social interactions should be designed in such a way that it can be worn exclusively by the user, and is able to monitor both sides of the interaction.

In two independent experiments [123] and [124], researchers developed a social feedback device that provides intervention when a person with visual impairment starts to rock their body in a stereotypical manner. [123] designed a device that consisted of a metal box with a mercury level switch that is worn on the user's head. A tone generator was activated when the user bent his upper body beyond a certain tilt angle. The authors tested the device on a congenitally blind individual who exhibited extreme body rocking and they conclude that the use of any assistive technology is useful only temporarily while the device is in use. They observe that the body rocking behavior returned to baseline levels as soon as the device was removed. Since the time of this experiment, behavioral psychology studies have explored short term feedback for rehabilitation [74], and these studies support the above observation that short-term feedback can even be detrimental to rehabilitation, and subject's case invariably worsens. Further, Jindal-snape emphasizes the need for long-term feedback that should target a slow, yet sustained rehabilitation to overcome stereotypic behaviors. Unfortunately, due to the prohibitively large design of the device developed by these researchers, it was impossible to have the individual wear the device over long

durations.

Observation 4: Assistive technology designed towards social assistance and behavioral rehabilitation should be used over long durations in such a way that the feedback is slowly tapered off over a significantly longer duration of time.

In [124] researchers used a ‘Drive Alert’ (i.e. a driver alerting system that audibly signals drivers when they start to fall asleep and their head droops forward) to detect body rocking, and to provide audio feedback to a congenitally blind 21 year old student. The researchers found that they were able to control body rocking effectively, but the device could not differentiate between body rocks and some other functional body movements. This device, primarily built to sense drooping in drivers, provides no opportunity to differentiate between a body rock and a functional droop. The resulting large number of false alarms discouraged the user from employing the device.

Observation 5: Assistive technology designed to facilitate behavioral rehabilitation should be effective in discriminating social stereotypic mannerisms from other functional movements to keep the motivation of device use high.

3.1.1 Design principles for social assistive and rehabilitative devices

Based on the above observations, a device developed to facilitate the social interactions of people with sensory or cognitive disabilities might do so by, (a) detecting social cues during social interactions and delivering that information to the user in real time to enable empathy, or (b) detecting the user’s own stereotypic behaviors during social interactions and communicating that information to the user in real time to provide social feedback. The first device would be classified as an assistive technology, while the second might be classified as a rehabilitative technology. Ideally, such a device would be based on the following design principles:

Design principle 1: The device should be portable and wearable so that it can be used in any social situation, and without any restriction on the user’s everyday life.

Design principle 2: The device should employ sensors and personal signaling devices that are unobtrusive, and do not become a social distraction.

Design principle 3: The device should include sensors that can detect the social mannerisms of both the user and other people with whom the user might communicate.

Design principle 4: The device should be comfortable enough to be worn repeatedly for extended periods of time, to allow it to be used effectively for rehabilitation.

Design principle 5: The device should be able to reliably distinguish between the user’s problematic stereotypic mannerisms and normal functional movements, to ensure that it will be worn long enough to achieve rehabilitation.

The remainder of this chapter looks at incorporating these design principles into a social interaction assistant targeted at enriching social situational awareness for individuals who are blind and visually impaired. The next section reviews relevant research aimed at identifying non-verbal cues with computer vision techniques.

3.2 Related Work in Computer Vision Research towards Sensing Factors that Contribute to the Overall Non-verbal Communication Picture

One of the goals of computer vision research is to develop pattern recognition and machine learning techniques that emulate the abilities of the human visual system. Computer vision research has matured over the past two decades to a point where many of the capabilities of human vision are now being emulated.

Table 3.1: Related Work in Automatic Detection of Interaction Environment Cues

Interaction Environment				
	Scene Change Detection	Background Modeling	Face and Object Detection	Environment Analysis
Proxemics		[125]	[126] [127]	
Objects in the scene	[128]	[129] [130]	[126]	
Natural vs manmade environment	[128]			[131]

Table 3.2: Related Work in Automatic Detection of 1) Physical Characteristics of the Communicator, and 2) Behavior of the Communicator

	Person Recognition	Clothing Recognition	Body Part Segmentation	Facial Feature Segmentation	Gender Race Recognition	Facial Motion Analysis	Body Motion Analysis	Eye Detection	Eye Tracking
Physical Characteristics of the Communicator									
Race and Body Color			[132] [133] [134]		[135]				
Body Shape	[136] [137]	[138] [139]	[140] [141] [132] [133] [134] [142]				[135] [137]		
Body Decoration	[143]								
Facial Hair				[144]					
Eye Glasses				[145]				[146]	
Clothing		[138] [139]							
Hair			[140] [147]						
Age	[148]								
Gender	[136]				[135]		[149]		
Identity	[150] [138] [151] [152] [153]					[154]			
Behavior of the Communicator									
Description of facial features				[155] [154]					
Body Mannerisms			[156] [157] [158]		[135]		[159] [160] [161] [162]		
Eye Gestures								[146]	[163] [164]
Gaze						[165]		[166] [167]	[165] [168]
Expressions and Emotions	[153]			[145] [155]		[116] [144] [145] [169] [170]	[171] [164] [162]	[172]	[163]
Personality	[136]		[157]		[135]		[161]		
Posture	[136]		[134]	[155]			[160] [144]		

The rows in Table 3.1 and 3.2 represents the various non-verbal cues listed in Section 1.3.1 through 1.3.3, while the columns represent the computer vision techniques being used in state-of-the-art computer vision research to detect those non-verbal cues. For example, Table 3.2 shows that [136] addressed posture by recognizing who the person is, and using that information to interpret the data accordingly. In contrast [134] used body part segmentation, [155] approached it through facial feature segmentation, and [160] and [144] evaluated posture by analyzing the motion patterns of the limbs.

3.3 Requirements Analysis for a Social Assistive Technology: Evidence Aggregation

As a part of the evidence-based model for developing the social interaction assistant, after identifying important design principles for assistive technologies (See Section 3.1), this section focuses on collecting evidence from the user community about the kinds of information that they felt would enrich their social interactions. The goal in doing this was to focus the development of the assistive technology human factor issues that would provide perceived quality of life improvements. In order to identify unmet needs of the visually impaired community, two focus groups were conducted¹. These groups consisted primarily of people who are blind, as well as disability specialists and parents of students with visual impairment and blindness where conducted.

The members of these focus groups who were blind or visually impaired were encouraged to speak freely about their challenges, in coping with daily living. During these focus groups, the participants agreed on many important problems. However, the problem of engaging freely with their sighted counterparts was highlighted as a particularly

¹ In order to understand the assistive technology requirements of people who are blind, we conducted two focus group studies (one in Tempe, Arizona USA - 9 participants, and another in Tucson, Arizona USA - 11 participants) which included:

1. Students and adult professionals who are blind,
2. Parents of individuals who are blind
3. Professionals who work in the area of blindness and visual impairments.

There was unanimous agreement among participants that a technology that would help people with visual impairment to identify people, or to hear a description of them would significantly enhance their social life.

important problem that was not being adequately addressed by assistive technologies².

While various other examples were cited by focus group participants, the inability to access non-verbal cues were considered to be of highest priority. These discussions produced the following list of 8 social needs often experienced by people with visual impairments.

1. Knowing how many people are standing in front you, and where each person is standing.
2. Knowing where a person is directing his/her attention.
3. Knowing the identities of the people standing in front of you.
4. Knowing something about the appearance of the people standing in front of you.
5. Knowing whether the physical appearance of a person who you know has changed since the last time you encountered him/her.
6. Knowing the facial expressions of the person standing in front of you.
7. Knowing the hand gestures and body motions of the person standing in front of you.
8. Knowing whether your personal mannerisms do not fit the behavioral norms and expectations of the sighted people with whom you will be interacting.

This list of 8 needs can be reduced to 2 basic categories of needs with regards to social interactions and social situational awareness were identified:

1. The need for access to the non-verbal cues of others during social interactions
2. The need for feedback about how one is perceived by others during social interactions.

²The following quotes are examples given by the focus group participants:

- “It would be nice to walk into a room and immediately get to know who are all in front of me before they start a conversation”.
- One young man said, “It would be great to walk into a bar and identify beautiful women”.

These two categories of needs were found to agree with the results of studies conducted by Jindal-Snape [74] with children who were visually impaired. She identified these two needs as *Social Learning* and *Social Feedback*, respectively.

To further understand the relative importance of non-verbal communication primitives, an online survey was carried out. The online survey consisted of eight statements that corresponded to the previously identified list of needs. Respondents indicated the level of their agreement with each statement using a Five-point Likert scale, the metrics being (1) Strongly disagree, (2) Disagree, (3) Neutral, (4) Agree, and (5) Strongly agree. The questions in the survey were used to infer the importance of non-verbal cues for social learning and social feedback. The survey was anonymously completed by 28 people, of whom 16 were blind, 9 had low vision, and 3 were sighted specialists in the area of visual impairment and vocational training.

3.3.1 Results from the Online Survey

3.3.1.1 Average Response

Table 3.3 shows the eight statements about social interactions that were evaluated with individuals who are blind and visually impaired. The results are sorted by descending importance, as indicated by the survey respondents (the question numbers correspond to the need listed in the previous section). The mean score is the average of the respondents on the 5 point scale that was used to capture the opinions. A score closer to 5 implies that the respondents strongly agree with a certain question and that they consider inaccessibility to that particular non-verbal cue to be important deterrent to their social interactions. On the other hand, a score closer to 1 represents the respondent did not consider the access to a specific non-verbal cue to be important during their social interactions.

3.3.1.2 Response on Individual Questions

Figure 3.1 shows the histogram of responses for the 8 Questions that were asked as part of the survey. Each subplot refers to a single question and shows the number of times users responded to that particular question with answers from 1 to 5 on the Lickert Scale. Each

Table 3.3: Average Likert Score on the 8 statements obtained through an Online Survey.

Question No.	Statements	Mean Score
8.	I would like to know if any of my personal mannerisms might interfere with my social interactions with others.	4.5
6.	I would like to know what facial expressions others are displaying while I am interacting with them.	4.4
3.	When I am standing in a group of people, I would like to know the names of the people around me.	4.3
7.	I would like to know what gestures or other body motions people are using while I am interacting with them.	4.2
1.	When I am standing in a group of people, I would like to know how many people there are, and where each person is.	4.1
2.	When I am standing in a group of people, I would like to know which way each person is facing, and which way they are looking.	4.0
5.	I would like to know if the appearance of others has changed (such as the addition of glasses or a new hair-do) since I last saw them.	3.5
4.	When I am communicating with other people, I would like to know what others look like.	3.4

histogram adds up to a total of 28 that corresponds to the 28 participants that took part in the online survey.

3.3.1.3 Response Ratio

Figure 3.2 shows the number of times the respondents chose to answer the 8 questions with their agreement or disagreement. The y-axis has been normalized to 100 points. The graph shows that respondents chose to answer the most by agreeing (Likert Scale 4) with the 8 questions. Followed closely behind was the strong agreement (Likert Scale 5) with the questions asked in the survey. The respondents chose to answer the least through strong disagreement (Likert Scale 1) to what was asked in the survey.

As described earlier, the 8 questions corresponding to the social needs of the indi-

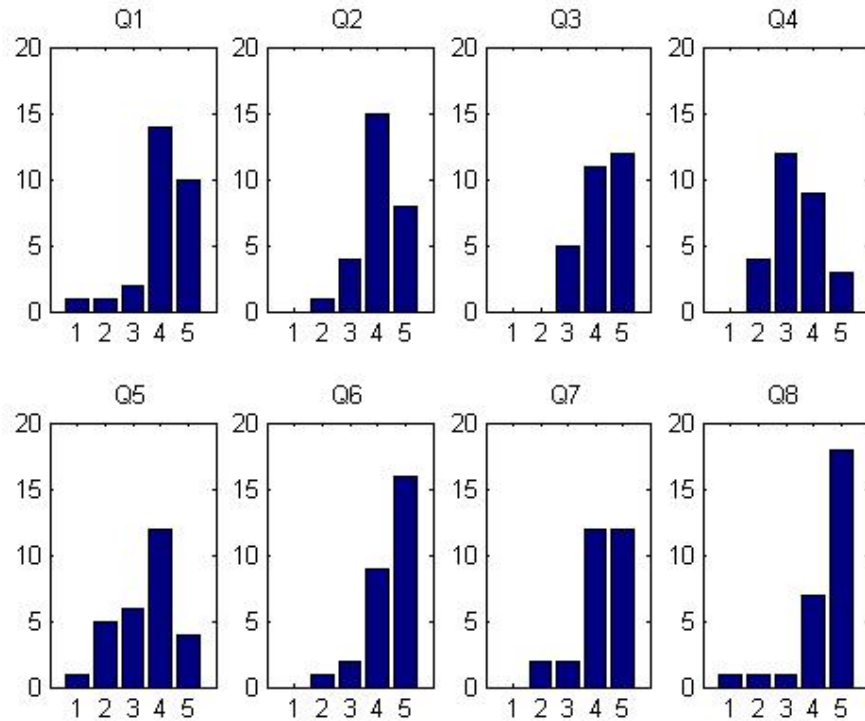


Figure 3.1: Histogram of Responses grouped by Questions

viduals were identified from the focus group survey that was conducted. Thus, the questions presented in the online survey questions were biased towards the needs of everyday social interactions of individuals who are blind and visually impaired. Thus, the implicit assumption while preparing this survey itself is that most of these items have been identified as being important and that only a priority scale needs to be extracted. This implicit assumption is immediately brought out by looking at the frequency with which the respondents answer with their agreement (Likert Scale 4) and strong agreement (Likert Scale 5).

3.3.1.4 Rank Average Importance Map for Various Non-verbal Cues

As explained earlier and as can be seen from Figure 3.2, the questionnaires were biased and the frequency of the responses is not Gaussian. This bias implies that using sample mean of the Lickert Scale responses will immediately show the same bias. This is due to the Gaussian iid assumption that is made while extracting the mean for the answers. In order to

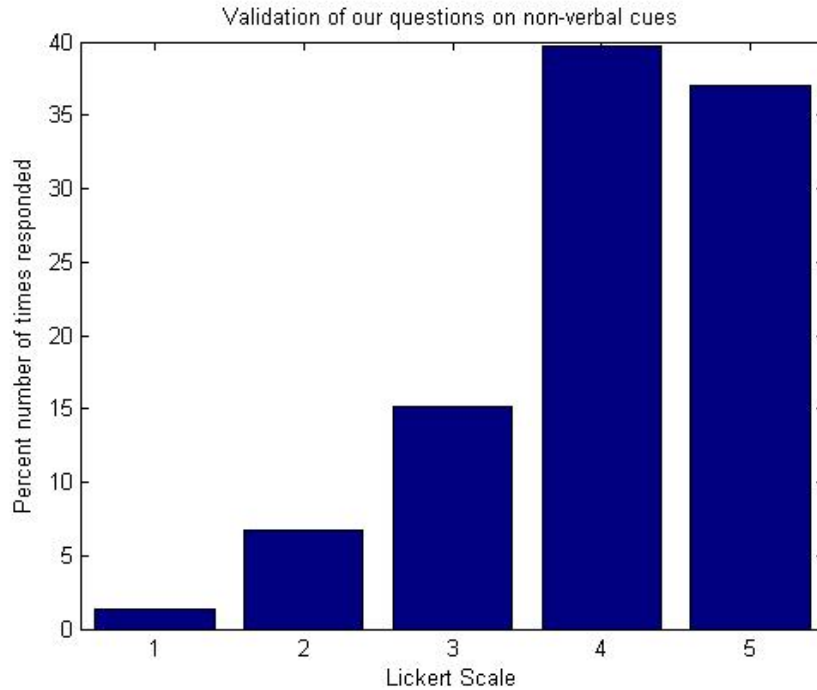


Figure 3.2: Response Ratio

overcome this non-Gaussianity, we resort to non-parametric mean for the responses. Rank average of the responses is estimated instead of the typical mean of the responses for each of the question. Please see Appendix A for the algorithm to determine the Rank Average. Since no assumptions on the distribution of the response are made, unlike the mean, the rank average gives a non-parametric method for comparing the responses of the individuals. The ranks can be either assigned ascending or descending with respect to the responses, i.e. rank 1 could mean all responses that were answered with strongly disagree (numeral 1), or rank 1 could mean all responses that were answered with strongly agree (numeral 5).

In the Figure 3.3, we have assigned rank 1 to strongly disagree. This is for the sake of visual convenience. Thus, higher the average rank, higher is that group's response from the respondents. Comparing Figure 3.3 to Table ??, it can be seen that the same ordering of priority can be seen through mean and rank average. But the mean tends to show very little variation between responses due to the bias that is present in the questions. On the other hand the rank average provides a good comparison scale.

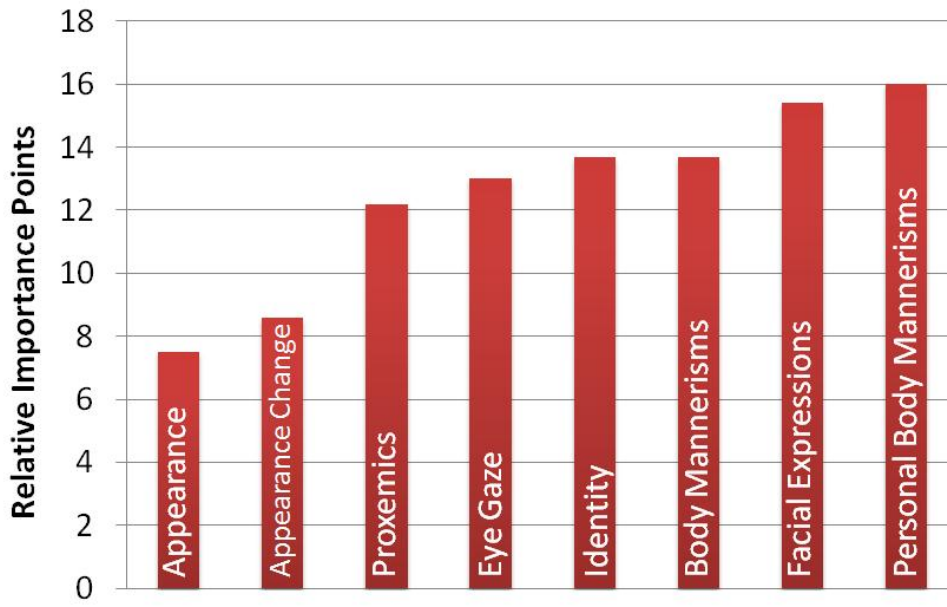


Figure 3.3: Rank average of the 8 questions

3.4 Evidence-based Model for the Proposed Social Interaction Assistant

The important observations from the above results include,

- Respondents are highly concerned about how their body mannerisms are perceived by their sighted peers (based on the response to Question 8 on the survey).
- Facial expressions form the most important visual non-verbal cue that individuals who are blind or visually impaired feel they do not have access to (based on Question 6 on the survey). This correlates with the studies into non-verbal communication that highlights the importance of facial mannerisms and gestures, which are mostly visual in their decoding.
- Followed by facial expressions, body mannerisms seem to be of higher importance for individuals who are blind and visually impaired (based on Question 3 of the survey).
- The responses to questions 7, 1 and 2 suggest that respondents would like to know the identities of the people with whom they are communicating, relative location of these

people and whether their attentions are focused on the respondent. This corresponds to knowing the position of their interaction partners when they are involved in a bilateral or group communication. People tend to move around, especially when they are standing, causing people who are blind to lose their bearing on where people were standing. This can result in individuals addressing an empty space assuming that someone was standing there based on their memory.

- The responses to questions 4 and 5 indicate that there was a wide variation in respondents' interest in (4) knowing the physical appearance of people with whom they are communicating and (5) knowing about changes in the physical appearance of people with whom they are communicating (See Figure 3.1. Many respondents indicated moderate, little, or no interest in either of these areas.

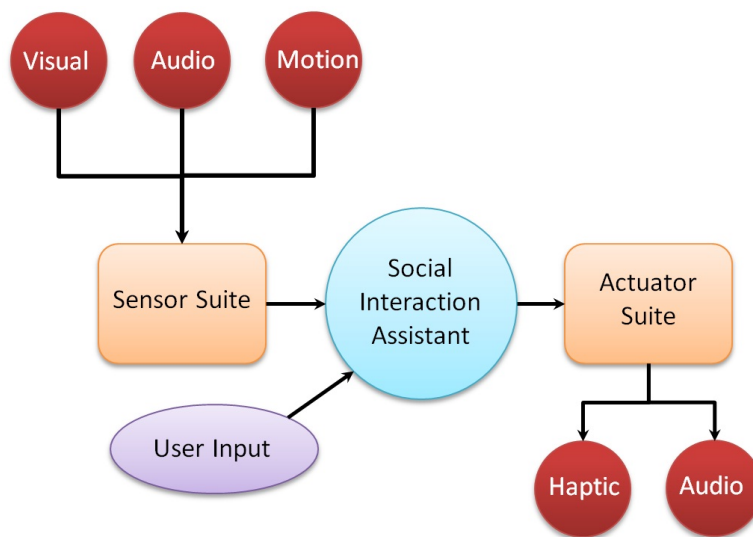


Figure 3.4: Social Interaction Assistant System Architecture.

3.5 System Architecture for a Social Interaction Assistant: An Evidence-based Assessment of Requirements

From the observations made on the development of a social interaction assistant, we conclude that the device should be portable, wearable and inconspicuous, and be able to provide access to visual cues, while enabling social feedback. To this end, we propose a system that

consists essentially of a set of body worn sensor suite capable of extracting important social signals from the environment, while at the same time using a suite of actuators that communicate extracted data back to the user without overloading their sensory system. Figure 3.4 shows the system architecture of the social interaction assistant, the realization of which is delayed until the components are described in the later chapters. Note the use of specific user input channel into the social interaction assistant along with the sensor suite. This provides unique opportunity for the user to control some of the basic functionalities of the system itself. The details of these user chosen functionalities will become apparent in the later chapters.

3.6 Organization of the Later Chapters

Following the organization of the dissertation presented in Table 1.5, Table 3.4 presents a detailed view of the chapters 4 through 10 that highlights the various social signal sensing and delivery within the context of developing a social mediation technology for helping people who are blind and visually impaired.

Table 3.4: Organization of the dissertation.

Chapter 4	Discusses the importance of providing social feedback to individuals who are disabled about their body mannerisms of how they are contributing to the social exchange that they are having with their sighted counterpart.
Chapter 5	Discusses the importance of facial expressions in social interactions and highlights technologies developed towards enabling access to dyadic interaction cues for users who maybe blind or visually impaired.
Chapter 6	Discusses and demonstrates technologies that can allow people who are blind and visually impaired to receive facial expressions of their interaction partners, with limited or no sensory loading.
Chapter 7	Discusses technologies that enable identification of interaction partners through facial biometrics. While developing this technology, the algorithms take advantage of the social context where they are being used and identifies people through person-specific facial landmarks, similar to human recognition of other individuals.
Chapter 8	Discusses technologies that could aid people who are blind and visually impaired to gain access to the social scene (number of people in front of them, their location, their identities and relative distance) in front of them.
Chapter 9	Discusses technologies that can track individuals who are in the social scene in front of the user who is blind and visually impaired. Tracking the individuals is important to ensure that the user's are provided with accurate information on the social scene structure.
Chapter 10	Discusses technologies that can deliver the social scene data, extracted as discussed in Chapter 8 & 9, while making sure that the user's do not face sensory overload.

Chapter 4

ENRICHING SOCIAL FEEDBACK: SENSING STEREOTYPIC BEHAVIORS

From the evidence-based modeling of the needs of the social interaction assistant, it was clear that people who are blind and visually impaired required social feedback on their body mannerisms more than any other social information. This is understandable as they don't receive any form of self-report on their behaviors if their sighted counterparts resort to unconscious non-verbal reactions that are visual in nature. It has been noticed that people who are blind become very conscious of this inability to receive social feedback and they resort to rather timid expression of body or face based gestures. While this prevents them from unconsciously displaying behaviors that may be considered atypical for the social setting, this also results in them appearing too rigid and asocial. That is, there exists a fine balance of bodily expressions which define the boundaries of cultural norms, which sensorially able people can learn through observation, but is inaccessible to people who are blind. In this chapter we discuss a variant of this above problem, residual stereotypic body rocking in adults who are blind and visually impaired. The problem is delineated in light of enriching social situational awareness and a solution is discussed along with a means of alleviating the condition.

4.1 Stereotypic Body Behaviors

Stereotypic behavior refers to any mannerism or utterance that is repetitive and non functional in nature [173] [174]. Stereotypy occurs in a large portion of the general population in various forms, although aggressive forms have been associated with sensory and cognitive disabilities in individuals. For example, individuals who are blind have the tendency to develop body rocking, head weaving, head drooping, and eye poking [175], while, individuals who are deaf have a propensity to develop various repetitive vocal behaviors [176][177]. Cognitive disabilities (both acquired, like brain injury and congenital, like autism and mental retardation) are associated with stereotypic behaviors like body rocking, hand flapping, jumping, and marching in place [178].

Though harmless by itself, Stereotypy can become a hindrance to social interactions and social acceptance [179]. Reference [124] introduces a 21 year old congenitally blind student who has an extreme case of body rocking (both while sitting and standing) that has become an obstruction to his career and an independent vocational evaluation states that a reduction in the student's body rocking was absolutely necessary for any form of employment. Stereotypy is a concerning problem in children, for whom peer acceptance is very important for their healthy growth and development of good social skills [180]. Children with stereotypic behaviors become victims of teasing thereby leading to social isolation, bullying and social segregation leading to negative self esteem. Aggravating these problems, social segregation and isolation have long term psychologically effects on the individual rendering an overall poor social skill set. Studies have shown that poor social skills are a leading cause for psychological problems such as depression, loneliness, and social anxiety [72].

Stereotypy, like any other human behavior, is very person specific. But socio-behavioral studies have shown that there are commonalities in these behaviors and there are broad classifications that can be identified in stereotypy prevalent in the general population. Eichel [181] introduces taxonomy for mannerisms that people with blindness and visual impairment tend to display. He identifies that body rocking appears on top of the most commonly seen behavior stereotype. A review of literature [182] further supports the claim that body rocking and head related mannerisms, including head weaving and drooping, are distinctive behaviors exhibited by individuals who have sensory or cognitive disabilities. For example, [175] discusses the case of a blind student who has developed extreme body rocking stereotype. The student bends in a 30 degree arc when he is sitting and when standing, places a foot well ahead of him and bends forward in an even greater arc. Such stereotypes can hinder the interactions of these individuals with friends and family, eventually leading to isolation and social inadequacy in their personal and professional life.

4.2 Focus of the chapter

Having identified stereotypic body behaviors to be an important deterrent in social acceptance of individuals with cognitive or physical disabilities, we focus on the possibility of building a rehabilitative and/or assistive technology towards providing feedback to individuals about their stereotypic body behavior. Specifically, we focus on body rocking, as it tops the list of most widely seen stereotypic behaviors. To this end, our research aims to answer three important questions:

1. Is there any evidence of individuals responding to rehabilitation for reducing stereotypic body rocking behavior?
2. If yes, what is the state-of-the-art technology available to detect and notify individuals of their rocking behavior?
3. Is it possible to build a device that detects body rocking condition and how well can it distinguish body rocking from other functional activities of daily living?

We answer the first question by looking into the immense literature available in behavioral psychology which has been studying behaviors in humans and their response to rehabilitation and assistance. To answer the second and third questions, we focus our attention towards wearable computing solutions that have gained a lot of momentum in the recent past. In specific, we develop an argument for an inclusive framework that uses state-of-the-art motion sensors with effective learning algorithms for detecting stereotype body rocking. As mentioned above, body rocking seems to be the most widely seen stereotype behavior and we use it as a basis for our argument that current level of technology can provide immense opportunities for developing rehabilitative and assistive technology solutions for reducing or controlling stereotypic behaviors.

4.3 Background and Related Work

4.3.1 *Foundations for social rehabilitation of behavioral stereotypes*

For over three decades, researchers in behavioral psychology have been publishing case studies on individuals who exhibit stereotypic body rocking. Most of these studies have targeted at reducing or controlling stereotypic body rocking. The methodologies used by these researchers, though varying in nature, can be broadly classified into two important categories.

4.3.1.1 *Intervention*

: Intervention relates to any form of feedback provided to an individual at the moment of exhibiting stereotype behaviors. Researchers have attempted to reduce body rocking by providing audio and/or tactual intervention whenever an individual started to rock. They have tried aversive punishment as well as less restrictive positive feedback in such situations. Felps and Devlin [124] issued an annoying tone in the ears of the subject while [182] used a recording of stone scratching on blackboard as the feedback tone whenever the individual started rocking. Both reported that the subjects responded well to the intervention. In contrast, [183], [183] and [184] have used verbal praise, physical guidance, verbal reprimands, and brief time-outs as intervention tools. Most of these researches have shown that intervention has worked in reducing and controlling body rocking without the use of aversive techniques. Aversive or not, these techniques validate a claim that it is possible to control or reduce body rocking (or any other stereotypic body mannerism) through feedback.

4.3.1.2 *Self Monitoring*

: In contrast to intervention, self-monitoring does not stop at intervening into the activities of the individual. It attempts to teach these individuals subtle cognitive skills to replace the current mannerism with more socially acceptable behavior, exercise, or medications. McAdam and O’Cleirigh [175] identifies that self monitoring is a very effective way of reducing the body rock behavior. They introduce the case of a congenitally blind individual

who is trained (with constant monitoring and positive feedback) to count the number of body rocks he goes through. Researchers noticed that the individual slowly waned off body rocking as he came to recognize and count his body's oscillatory movements. The research concludes that a well designed self monitoring program could benefit in reducing stereotypic body rocking. Shabani, Wilder and Flood [178] presents the case of a 12 year old child who was diagnosed with attention deficit hyperactivity disorder (ADHD) having an excessive body rocking and hand flapping stereotypy. The authors introduce an elaborate and positively rewarding self monitoring scheme that allows the child to improve on his behavior effectively. A follow-up with the child's teacher indicated that the social outlook of the child had improved over the course of rehabilitation and the case further reiterates ability to rehabilitate individuals with stereotypic behavior. Estevis and Koenig [185] introduces a cognitive approach to reducing body rocking on an 8 year old congenitally blind child through self monitoring. Teachers or family members would tap on the shoulders of the child when he started rocking, while the child was taught to recite his own monitoring script. The authors conclude that rocking can be significantly reduced through notification to the individual combined with self monitoring.

Supporting such case studies of behavioral mannerisms, psychologists have been studying intervention and feedback as an integral component of social development. Feedback can be defined as the provision of evaluative information to an individual with the aim of either maintaining present behavior or improving future behavior [186]. According to [187], feedback is critical to social development because after an individual receives information about his or her performance, he or she can make the necessary modifications to improve social skills. Most social skills develop during early years and in order for children to evaluate themselves accurately and to modify social skills, it is essential that children to be given feedback [73] [75], since without clear feedback, the children are unable to identify how their social behavior differs from others or is perceived by others in the environment [188]. Based on these studies there is enough evidence that feedback that offers intervention, possibly followed by a well planned self-monitoring program could benefit in reducing or controlling body rocking behavior.

4.3.2 *Need for Assistive or Rehabilitative Technology*

The feedback needed for intervention usually comes from people in and around these individuals who have stereotypic behavior. It has been observed that significant others in the environment often fail to give feedback, and even when they do, it is not meaningful or understandable to individuals who need rehabilitation - for example, in case of individuals who are blind or visually impaired, nodding one's head in reply to a question or gesturing [74] would be futile. Meaningful feedback is important, not only for social interaction, but for accurate self-evaluation by individuals. Most times people within the vicinity of individuals with needs fail to offer these crucial feedbacks. Many times, the individuals with needs feel guilty or obligated to ask for help from others in their environment. The ability to augment or replace this significant individual(s) in the environment with a reliable feedback mechanism is the aim and goal of all assistive technology solutions (In an independent online survey conducted by [189], the researchers found that people who are visually impaired would expressed the need for an assistive technology that would provide feedback on their own social mannerisms and offer a potential to improve their social outlook). Focusing on the development of such a technology that effectively detects body rocking and provides feedback to an individual is the goal of this paper. While we focus only on intervention through feedback, in the Future works section we highlight some ideas for extending the proposed framework into self-monitoring tools.

4.3.2.1 Past research into building assistive technology to detect body rocking

Transon [123] developed a head mounted switching device that would trigger a tone when an individual starts to rock. The device consisted of a metal box with a mercury level switch that detects any bending actions. The feedback was provided with a tone generator that was also located inside the metal box. The entire box was mounted on a strap that the user wears around his/her head such that the speaker aligns with the ears. The authors tested it on a congenitally blind individual who had severe case of body rocking and they conclude that the use of any assistive technology is useful only temporarily while the device is in

use. They state that the body rocking behavior returned to baseline levels as soon as the device was removed. Since the time of this experiment, behavioral psychology studies have explored short term feedback for rehabilitation [74] and these studies support the above observation that short time feedback is most of the times detrimental to rehabilitation and subject's case invariably worsens. Unfortunately, due to the prohibitively large design of the device developed by these researchers, it was impossible to have the individual wear the device over long durations. Thus, any technology developed for behavioral rehabilitation should be small and researchers should target the use over long durations in such a way that the feedback is slowly tapered off over a significantly longer duration of time.

Similar to the previous experiment, [124] used a 'Drive Alert' (driver alerting system that monitors head droop) to detect body rocking and provide feedback to a congenitally blind 21 year old student. The research concludes that they were able to control body rocking effectively, but the device could not differentiate between body rocks from any other functional body movements. This device, primarily built to sense drooping in drivers provides no opportunity to differentiate between a body rock and a droop. Use of such devices could only be negative on the user as a large number of false alarms would only discourage an individual from using any assistive technology. Assessing these above technologies, we resort to two important design dimensions in every step of the building of our assistive device.

1. Size and placement of the device: We argue that any assistive device developed for the sake of improving social outlook of an individual should respect the appearance of a person in his/her social circle and should provide a solution that is discrete and non intrusive. We call this the Acceptance dimension.
2. Ability to discriminate rocking from other functional activities: False feedback even over a short period of time could be discouraging for an individual to continue using his/her assistive tool. It is imperative that the device be able to distinguish between the stereotypy from any other form functional activities effectively to keep the motivation of device use high. We call this the Motivation dimension.

The proposed methodology uses these two design dimensions while addressing the need of a new assistive technology.

4.4 Methodology

Recently, human activity detection and recognition using motion sensors have taken a front seat in technology and behavioral research. This is due to the availability of micro mechanized electronic systems (MEMS) that have started to implement complex mechanical systems at a micro scale on integrated circuit chips. These offers advantages like reliability, cheaper cost of production, smaller form factor and above all extremely precise measurement with least or no maintenance. One such sensor is the accelerometer that is capable of measuring the effect of gravity on three perpendicular axes. When mounted on any moving object, the opposing motion (opposing gravity) of the entity allows these sensors to measure the speed and direction of motion. Integrating the magnitude and orientation information over time it is possible to accurately measure the exact motion pattern of the moving entity. These accelerometers have been used by researchers to track motion activity in almost every joint of the human body [190]. Researchers have used single, double or triple orthogonal axis accelerometers to detect various activities of humans. They all follow the same underlying supervised learning architecture with difference in learning algorithm used. A simplified representation of the same is shown in Figure 4.1.

Five bi-axial accelerometers are used in [190], along with a decision tree classifier to detect and recognize 20 different activities of daily life. They report a recognition rate of over 85%. In [191], the authors evaluated different meta classifiers for recognizing seven lower body motion patterns from a single biaxial accelerometer data and reported the best performance for boosted Support Vector Machines (SVM) [192] with a subject independent accuracy of 64%. Since each dimension of the accelerometer data is similar to audio waveform, popular Hidden Markov Models [193] can be used to learn motion patterns. Reference [194] used HMM to learn the accelerometer data for specific tasks performed by participants and reports a recognition rate of over 90%. In [195], researchers have used two accelerometers placed on the arms of Kung-Fu practitioner and report a recognition accu-

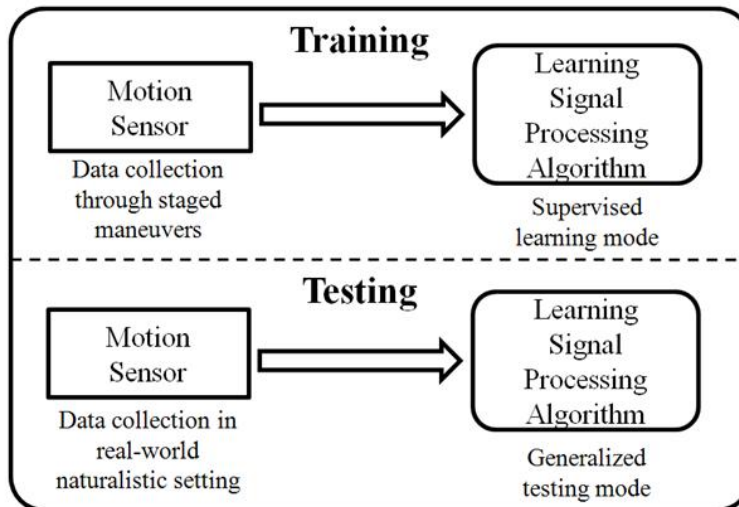


Figure 4.1: Training and testing phases of a typical learning framework found in literature.

racy of 3 Kung-Fu arm movements at 96.6%. Research work [196] demonstrates the use of accelerometer data to not only recognize activity, but also localize people within a building. Though the technique is rudimentary, the authors report a high accuracy in recognition of activities while localization still remains a research topic. [197] have demonstrated the use of accelerometers in not only monitoring movements, but also static posture of the human body. They report a recognition rate of 95% using four sensors placed on the chest, thigh, forearm and wrist of participants. Extending this work, [198] have demonstrated an assistive technology solution that uses low cost accelerometers on stroke patients and monitor their posture and walking patterns. Using this information, a feedback is provided to the patient to self-correct their posture and walking pattern.

Based on all these findings, we hypothesize that an accelerometer based motion detector should be capable of capturing body rocking data and should be able to discriminate between rocking and other functional activities. We specifically chose the motion sensor and learning algorithm based on previous work done at our institute with the detection of seven simple body activities [199]. Researchers analyzed the performance of discriminative classifiers like AdaBoost, Support Vector Machines and RLogReg for recognizing these seven different activities and concluded that AdaBoost classifier offered the best recognition

rate at 94%. Based on these results, in this paper, we extend the use of AdaBoost learning framework into body rock detection. We discuss the use of two AdaBoost classifiers - the classical AdaBoost [200] and the more recent Modest AdaBoost [201] for detecting and discriminating body rocking effectively. Our focus in the paper is directed towards understanding the generalization capabilities of the two AdaBoost learning models so that false positive rate is reduced while keeping the true positive rate high.

4.4.1 Motion Sensors - Design choice along the "Acceptance" Dimension

In order to keep the motion detector discrete, we have chosen state-of-the-art tri-axial accelerometer package, ZStar III [202], marketed by Freescale Semiconductor. The accelerometer is shown in the inset of Figure 4.2. The device (including a coin battery as a power source) is an inch in diameter and less than eighth of an inch in thickness thereby allowing an elegant integration into everyday clothing. Figure 4.2 shows the typical use of the accelerometer in the proposed application for detecting body rocking. The accelerometer has a very high sensitivity with protection against excessive g-force damage. The sensors wirelessly connect to a PDA and/or cell phone through IEEE 802.15.4 (ZigBee) wireless standards. The use of low power consumption electronics for both acceleration sensing and wireless communication allows this device to work for hours at length on a single coin battery. Further, the advanced sleep mode implementations allow the device to stay at nano watt power mode during non-operation. The proposed solution allows for prolonged use of the device to the effect of an assistive technology thereby maintaining a longer duration feedback based rehabilitation regimen.

processing element for the current study was a Windows Mobile Operating System based PDA running on a 400Mhz XScale processor. The software components (described in detail in Section III-B of the proposed solution were placed on the PDA that could be carried by a user without any extra load. The proposed assistive technology is an addition to the Social Interaction Assistant proposed in [203]. The software component implementation is generic to be ported to most modern cell phones that possess enough processing power, but is always under utilized for its capacity. The feedback (an audio tone) is currently

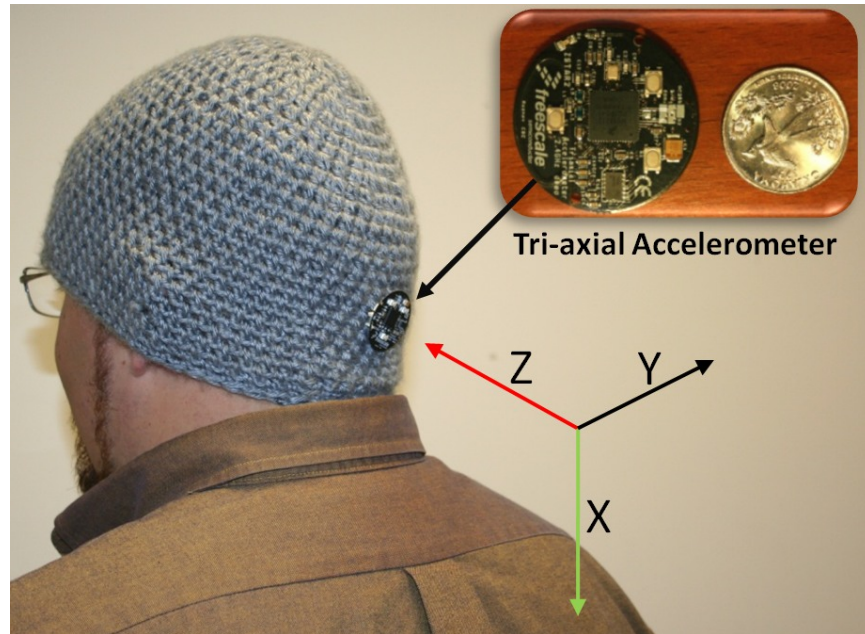


Figure 4.2: The proposed hardware for use in the detection of body rocking stereotypic behavior. The accelerometer, in comparison with a US quarter, is shown in the inset. The three axes marked in the image shows the orientation of the accelerometer as it is placed on the head.

being provided through a Bluetooth headset that is paired with the processing element. The choice of this feedback device was again based on the idea that Bluetooth headset has everyday acceptance among the masses and is no longer seen as a social distraction. In future, we plan to explore the use of delivery modalities that transcends the typical visual and audio medium. We intend to use haptic cues to inform the participant not only their rocking behavior but more complex self-monitoring routines that could allow the user to withdraw from the rocking behavior effectively.

Figure 4.3 shows a typical data stream collected from the accelerometer shown in Figure 2 during rocking and non-rocking functional behavior. The three data streams correspond to the three axis of the accelerometer each sampled at 100 Hz. It can be seen that the data stream under rocking conditions are visually distinguishable when compared to non-rocking functional movements. The following section highlights our choice of learning framework and features we extracted from these data stream in order to achieve reliable rocking and non-rocking discrimination.

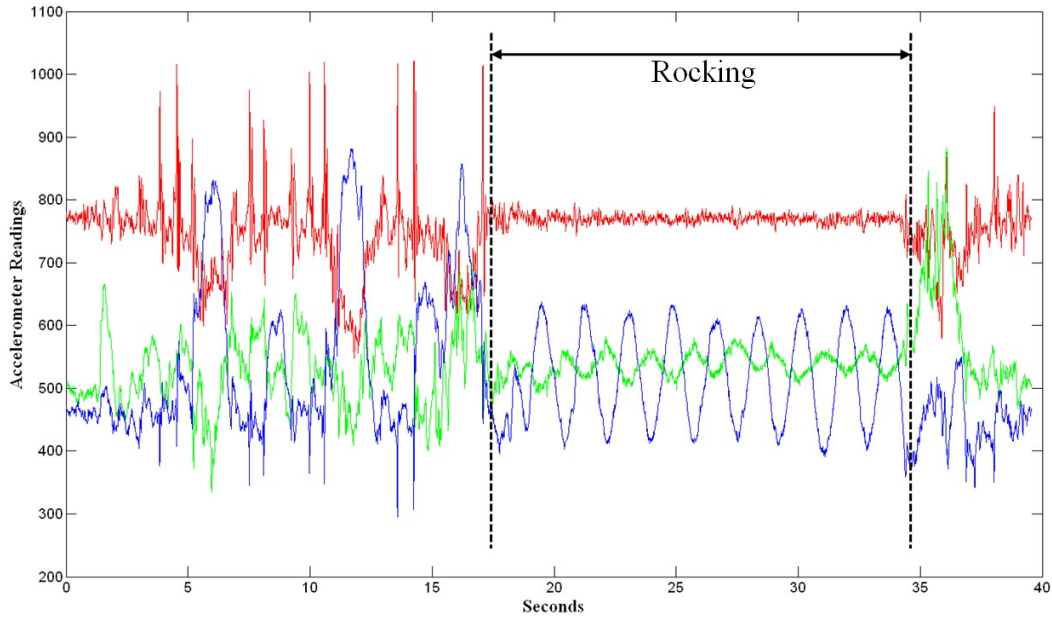


Figure 4.3: Data stream for the tri-axial accelerometer. The three streams correspond to the three axes. The figure shows non-rocking events followed by rocking and then followed by non-rocking.

4.4.2 *Extracting Body Rock Information from Motion Sensor Data - Design choice along the “Motivation” Dimension*

As mentioned before, the work presented in this paper builds on top of the work presented in [199] where the authors use two accelerometers placed one at the ankle and the other on the thigh to distinguish between simple activities like walking, running, standing etc. They proved the use of an aggregated AdaBoost classifier system that was built out of simple linear classifiers to achieve activity recognition. Unfortunately, the work does not provide any assessment on the generalization capabilities of their aggregate classifier. We extend their work into the problem of body rock detection using only one accelerometer placed on the back of the person’s head. Below, we discuss the various features that we extract from the accelerometer data and introduce the variant of AdaBoost that generalizes on its training set very well (termed Modest AdaBoost). We show results of our experiments and discuss our reasoning to believe how the new AdaBoost framework is able to generalize on body rocking data when compared to classical AdaBoost used by [199].

4.4.2.1 Features:

Since we are using a tri-axial accelerometer, we obtain three orthogonal axis data through rocking and non-rocking events. In order to capture the temporal variation in the acceleration data, we accumulate the input stream on each axis for a fixed duration T seconds and all features are extracted on this packet of acceleration data. As a part of the assessment, we determine the best packet length for the task of body rock detection. Further, successive packets are extracted with a fixed duration of overlap between them.

We chose five sets of features that were extracted on the three axes of accelerometer data. For the sake of clarity, we cluster these sets into two groups based on whether they were chosen due to popular use in the accelerometer data processing community or due to the author's insights into the body rocking data.

Group 1 - Popular features used by the motion analysis research community [190] [199]: We choose the following three sets each of which were applied on all three axes of acceleration data, henceforth referred to as x, y, z axis data.

1. Mean of x, y, z data over the duration of packet.
2. Variance of x, y, z data over the duration of packet.
3. Correlation between the three axes (x-y, y-z and z-x) over the duration of packet.

Group 2 - Authors insights into body rocking data: Inspecting the accelerometer data shown in Figure 3, it can be seen that the Z axis changes from random signal pattern to more of a sinusoidal pattern when the individual's behavior transitions from non-rocking to rocking. Thus we choose two sets of features which we hope would capture this non-sinusoid to sinusoid transition between events. These features include

4. The first order differential power on all three axes - Sinusoidal signals change gradually over time such that the averaged sum square energy in the temporal first order differential

of the signal should be less when compared to a random signal where the first order differential can have very high variations and hence higher power.

5. Fourier Transform variance and kurtosis on the Z-axis only - An effective way to capture power distribution of a signal into sinusoids is by using Fourier Transform. We hypothesize that the non-sinusoid to sinusoid transitions can be captured by quantitatively measuring the power spread spectrum of the Z-axis accelerometer data. We model the power spread to be a Gaussian and extract the variance and kurtosis (peaking) of the spread to determine if there is rocking or not.

Table 4.1: Features for Body Rock Detection: Group 1

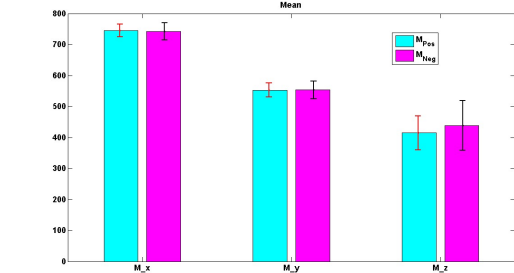
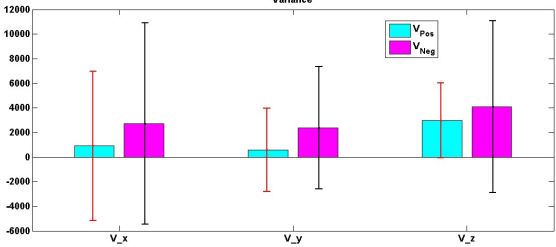
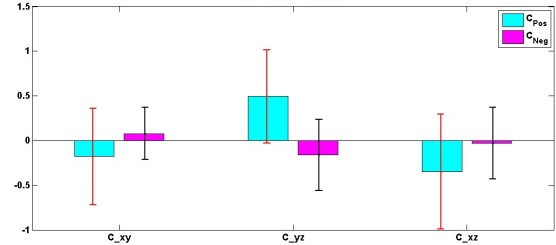
Group 1															
<p>Set 1</p> <p>Definition: Mean on the temporal dimension. Axes affected: x, y, z Number of contributing features: 3 Feature Identification Numbers: 1, 2, 3</p>	$M_x = \frac{1}{N} \sum_{i=1}^N x_i$	<p>1. M_x</p> <p>2. M_y</p> <p>3. M_z</p>	 <table border="1"> <caption>Mean Feature Values</caption> <thead> <tr> <th>Axis</th> <th>M_{Pos}</th> <th>M_{Neg}</th> </tr> </thead> <tbody> <tr> <td>M_x</td> <td>~750</td> <td>~750</td> </tr> <tr> <td>M_y</td> <td>~550</td> <td>~550</td> </tr> <tr> <td>M_z</td> <td>~420</td> <td>~450</td> </tr> </tbody> </table>	Axis	M_{Pos}	M_{Neg}	M_x	~750	~750	M_y	~550	~550	M_z	~420	~450
Axis	M_{Pos}	M_{Neg}													
M_x	~750	~750													
M_y	~550	~550													
M_z	~420	~450													
<p>Set 2</p> <p>Definition: Variance on the temporal dimension. Axes affected: x, y, z Number of contributing features: 3 Feature Identification Numbers: 4, 5, 6</p>	$V_x = \frac{1}{N-1} \sum_{i=1}^N (x_i - M_x)^2$	<p>4. V_x</p> <p>5. V_y</p> <p>6. V_z</p>	 <table border="1"> <caption>Variance Feature Values</caption> <thead> <tr> <th>Axis</th> <th>V_{Pos}</th> <th>V_{Neg}</th> </tr> </thead> <tbody> <tr> <td>V_x</td> <td>~1000</td> <td>~3000</td> </tr> <tr> <td>V_y</td> <td>~1000</td> <td>~2500</td> </tr> <tr> <td>V_z</td> <td>~3000</td> <td>~4000</td> </tr> </tbody> </table>	Axis	V_{Pos}	V_{Neg}	V_x	~1000	~3000	V_y	~1000	~2500	V_z	~3000	~4000
Axis	V_{Pos}	V_{Neg}													
V_x	~1000	~3000													
V_y	~1000	~2500													
V_z	~3000	~4000													
<p>Set 3</p> <p>Definition: Cross Correlation between axes. Axes affected: x, y, z Number of contributing features: 3 Feature Identification Numbers: 7, 8, 9</p>	$C_{xy} = \frac{1}{N-1} \sum_{i=1}^N (x_i - M_x)(y_i - M_y)$	<p>7. C_{xy}</p> <p>8. C_{yz}</p> <p>9. C_{xz}</p>	 <table border="1"> <caption>Correlation Coefficient Feature Values</caption> <thead> <tr> <th>Pair</th> <th>C_{Pos}</th> <th>C_{Neg}</th> </tr> </thead> <tbody> <tr> <td>C_{xy}</td> <td>~-0.2</td> <td>~0.1</td> </tr> <tr> <td>C_{yz}</td> <td>~0.5</td> <td>~-0.2</td> </tr> <tr> <td>C_{xz}</td> <td>~-0.2</td> <td>~0.0</td> </tr> </tbody> </table>	Pair	C_{Pos}	C_{Neg}	C_{xy}	~-0.2	~0.1	C_{yz}	~0.5	~-0.2	C_{xz}	~-0.2	~0.0
Pair	C_{Pos}	C_{Neg}													
C_{xy}	~-0.2	~0.1													
C_{yz}	~0.5	~-0.2													
C_{xz}	~-0.2	~0.0													

Table 4.2: Features for Body Rock Detection: Group 2

Group 2																												
<p>Set 4</p> <p>Definition: First order differential power. Axes affected: x, y, z Number of contributing features: 3 Feature Identification Numbers: 10, 11, 12</p>	$D_x = \sqrt{\sum_{i=2}^N (x_i - x_{i-1})^2}$	<p>10. D_x</p> <p>11. D_y</p> <p>12. D_z</p>	<table border="1"> <caption>Approximate data for First order differential power</caption> <thead> <tr> <th>Feature</th> <th>Class</th> <th>Mean</th> <th>Variance</th> </tr> </thead> <tbody> <tr> <td rowspan="2">D_x</td> <td>Pos</td> <td>~85</td> <td>~100</td> </tr> <tr> <td>Neg</td> <td>~220</td> <td>~150</td> </tr> <tr> <td rowspan="2">D_y</td> <td>Pos</td> <td>~75</td> <td>~80</td> </tr> <tr> <td>Neg</td> <td>~135</td> <td>~100</td> </tr> <tr> <td rowspan="2">D_z</td> <td>Pos</td> <td>~70</td> <td>~80</td> </tr> <tr> <td>Neg</td> <td>~140</td> <td>~100</td> </tr> </tbody> </table>	Feature	Class	Mean	Variance	D_x	Pos	~85	~100	Neg	~220	~150	D_y	Pos	~75	~80	Neg	~135	~100	D_z	Pos	~70	~80	Neg	~140	~100
Feature	Class	Mean	Variance																									
D_x	Pos	~85	~100																									
	Neg	~220	~150																									
D_y	Pos	~75	~80																									
	Neg	~135	~100																									
D_z	Pos	~70	~80																									
	Neg	~140	~100																									
<p>Set 5</p> <p>Definition: Gaussian fit power spread spectrum - Variance and Kurtosis. Axes affected: z Number of contributing features: 2 Feature Identification Numbers: 13, 14</p>	<p>If, $Freq_k = \{-\frac{\gamma}{2}, \dots, 0, \dots, \frac{\gamma}{2}\}$</p> <p>$\gamma$ is the sampling Frequency</p> $X_k = \sum_{i=1}^N x_i e^{-\frac{2\pi i}{N} kn}, k = \{1, \dots, N\}$ <p>then</p> <p>FFT Variance: $F_v = \sum_{i=1}^N X_k(Freq_k)^2$</p> <p>FFT Kurtosis: $F_k = \sum_{i=1}^N X_k(Freq_k)^4$</p>	<p>13. F_{v_z}</p> <p>14. F_{k_z}</p>	<table border="1"> <caption>Approximate data for FFT variance and kurtosis</caption> <thead> <tr> <th>Feature</th> <th>Class</th> <th>Mean</th> <th>Variance</th> </tr> </thead> <tbody> <tr> <td rowspan="2">F_{v_z}</td> <td>Pos</td> <td>~270</td> <td>~100</td> </tr> <tr> <td>Neg</td> <td>~260</td> <td>~100</td> </tr> <tr> <td rowspan="2">F_{k_z}</td> <td>Pos</td> <td>~5.5</td> <td>~1.5</td> </tr> <tr> <td>Neg</td> <td>~5.0</td> <td>~1.5</td> </tr> </tbody> </table>	Feature	Class	Mean	Variance	F_{v_z}	Pos	~270	~100	Neg	~260	~100	F_{k_z}	Pos	~5.5	~1.5	Neg	~5.0	~1.5							
Feature	Class	Mean	Variance																									
F_{v_z}	Pos	~270	~100																									
	Neg	~260	~100																									
F_{k_z}	Pos	~5.5	~1.5																									
	Neg	~5.0	~1.5																									

70

The figures shown in the last column plots mean values of data from positive rocking samples and negative rocking samples as bars. The variance on the same is shown as vertical error lines around the mean. The lighter (blue when viewed in color) shaded bar are values from the positive class, whereas the darker (pink when viewed in color) bar are values from the negative class.

Thus, the features used in our study can be categorized as belonging to two groups with three sets in Group 1 and two sets in Group 2. Each set has varying number of features based on what parameter the set is extracting from the temporal accelerometer data. Based on the descriptions above, the entire feature set has a total of 14 features. We identify each of these by their respective Feature Identification Numbers. Table 1 shows the two groups and the different sets under the group with typical values of these features under rocking and non-rocking behavior.

4.4.2.2 Learning Algorithm:

As discussed in introduction of this section, we compare the performance of two AdaBoost learning frameworks to determine which one can generalize the best on the training data. The two algorithms are introduced briefly below. For further details, the reader is referred to appropriate references provided within the subsections.

(a) Classical AdaBoost Learning Framework: AdaBoost learns any classification problem by working with a set of weak classifiers. Weak classifiers are those classifiers that use simple decision steps to categorize data into one of two pools - positives or negatives (In all our experiments, we used a three level decision tree [204] as the simple classifier). AdaBoost proceeds by ranking the labeled training data as being simple to complex based on how many weak classifiers are needed to learn each of the examples. The process continues on an iterative manner until all the training examples are learnt or till the allowed number of learning cycles are exhausted. Let, X be the input to a learning algorithm, in our case the features extracted as explained in the previous step, and Y be the label of what class the data belongs to, in our case, $Y = \{1, -1\}$ implying rocking, non-rocking, respectively. Values at each dimension of input X can be considered to characterize the incoming data in some manner and the task of the learning algorithm is to learn these representational values of the input dimensions that allow the algorithm to distinguish between rocking and non-rocking. AdaBoost does this learning by using a large set of simple (weak) learners (or classifiers) that act on each of the dimension of the input data with the determined goal of distinguishing rocking from non-rocking. The final decision of the complete learning mod-

ule is a combined opinion of all the simple learners that make up the system. The beauty of AdaBoost implementation is that the human intervention into the learning process stops at identifying what simple (weak) learners to use and what feature pool to operate on. Selection of number of weak learners, selection of input dimension on which the weak learners have to act, and the confidence to place on the decision of each of the weak learner is all determined by the algorithm during the training phase. Once the algorithm is trained, the final learnt rocking/non-rocking classifier can be represented as

$$L(x) = \text{sign} \left[\sum_{i=1}^N w_i f_i(x) \right] \quad (4.1)$$

where, x : An instance of all possible rocking patterns X . L : The final learnt classifier that can distinguish input x as rocking or non-rocking. f : The simple (weak) learner. N : The total number of weak learners that make up the complete learner L . w : Weight associated with each weak learners output. This corresponds to the confidence placed in each weak learner by the Boosted system.

From a learning perspective, in each step of the iterative learning, the AdaBoost algorithm implements a greedy optimization to pick a set of weak learners that minimize exponential classification error of the picked simple classifiers as shown below

$$Error_k = \sum_{i=1}^M e^{-y_i \cdot L(x_i)} \quad (4.2)$$

where, y : Label of the input instance x M : Total number of examples in the training set k : Learning iteration number

Further, based on each iterative step, a distribution (D_m) is created over the training set examples to represent their complexity (difficulty to learn). For example, in a given iteration, an example that could be solved is assigned a lower distribution weight while, a sample that was not learnt in that iteration step is assigned a higher weight. The lower weight on the learnt example implies that this example will be stressed less in the next learning iteration while all other examples which could not be solved will become the focus

for picking new weak learners. Moving from one iteration to the next, all the weak learners from the past k iterations are added into a pool of selected weak learners leading up to the final classifier L .

(a) *Modest AdaBoost Learning Framework*: All learning algorithms, including AdaBoost suffer from the problem of over fitting or over learning. This is due to the fact that training sample sets of positives and negatives can never be representative of all the possible samples that the algorithm will face in its operational life span. Since the learning is limited to a restricted set of examples, there is always the problem of over fitting into this small set and thereby loosing the ability to generalize their learnt knowledge to all other possible examples. To this end, many alternatives have been proposed to AdaBoost that will allow the algorithm to generalize better. We introduce Modest AdaBoost [201] which was recently proposed towards better generalization capabilities and has been shown to be powerful on various machine learning datasets. Unlike the classic AdaBoost where the distribution penalizes only examples that are not learnt in the previous iteration, Modest AdaBoost penalizes for examples that are not learnt and also examples that are learnt very well (over fitting). This is done by projecting all the examples in the training pool on to four separate distributions,

1. $P_m^{(+1)} = P_{(D_m)}(y = +1 \cap L(x))$: Probability of the learner, as measured on D_m , predicting an input instance x correctly as being rocking when the label also represents it to be rocking.
2. $P_m^{(-1)} = P_{(D_m)}(y = -1 \cap L(x))$: Probability of the learner, as measured on D_m , predicting an input instance x correctly as being non-rocking when the label also represents it to be non-rocking.
3. $\bar{P}_m^{(+1)} = P_{(\bar{D}_m)}(y = +1 \cap L(x))$: Probability of the learner, as measured in the inverse distribution (\bar{D}_m), predicting an input instance x correctly as being rocking when the label also represents it to be rocking.

4. $\bar{P}_m^{(-1)} = P_{(\bar{D}_m)}(y = -1 \cap L(x))$: Probability of the learner, measured in the inverse distribution (\bar{D}_m), predicting an input instance x correctly as being rocking when the label also represents it to be rocking.

Conditions 1 and 2 penalize the classifier on examples that are not learnt during a training iteration, whereas 3 and 4 penalize examples that are already learnt in the previous iteration which was learnt again in the current iteration. Combining these four measures as

$$f_m = \left(P_m^{(+1)}(1 - \bar{P}_m^{(+1)}) - P_m^{(-1)}(1 - \bar{P}_m^{(-1)}) \right) (x) \quad (4.3)$$

provides a means for penalizing the learner for not classifying an example and also for over fitting an example. This provides a means for modest learning of the final combined classifier L . We hypothesize that the choice of a learning algorithm that generalizes well will provide the opportunity to allow better non-rocking detection thereby hopefully increasing discrimination ability for the assistive device. This would directly reflect upon the motivation of the user to get feedback only when he/she is rocking and not performing other functional activities.

4.5 Data Collection

Two separate data collections were carried out, one in a controlled setting while the other in a more uncontrolled naturalistic everyday research laboratory setting. The controlled setting data collection was used for training and lab testing the device, whereas the uncontrolled naturalistic setting was used to determine how well the learning algorithm was able to generalize when used for an extended period of time as an assistive tool.

4.5.1 *Controlled Data Collection:*

Data was collected on ten participants who did not have any known stereotype rocking behavior. The goal of the experiments was to collect data for training the system to differentiate rocking from non-rocking behavior. To this end, we devised three separate data

collection routines where the subjects were required to do rocking and non-rocking tasks as naturally as possible. The details of the routines are as follows:

4.5.1.1 Routine A: Rocking data

Participants were allowed to choose from a rocking chair or a stool or sitting on the ground, so they could rock as comfortably and naturally as possible. We found some cultural preferences to the way people choose to rock. The subjects were asked to rock for a total of 20 complete cycles.

4.5.1.2 Routine B: Non-rocking data

The participants were asked to do activities that did not involve rocking. They moved around the experimental setup reading posters, operating computers, interacting with everyday office equipments and included some functional body motions similar to rocking like, stooping down to pick up objects, rapidly bending down to pick up objects etc. Data was collected for a total of 30 seconds.

4.5.1.3 Routine C: Test data

Since rocking can happen at any given instance, we collected data where subjects did various activities and interspersed them randomly with rocking. The goal is to determine how fast and accurately our system can detect such rocking occurrences. In all of these data streams, rocking instances were manually identified and marked for the sake of ground truth. Figure 3 shows the combination of rocking and non-rocking activities by the participants. It can be noticed that there is a clear demarcation between the two activity zones.

4.5.2 *Uncontrolled Data Collection:*

The uncontrolled data was collected towards testing the generalization capabilities of the learnt system. To this end, the body rock detection system was worn by the primary author during everyday laboratory activities. Body rock detection was provided as a feedback through a pair of headphones in the form of an audio beep. Five trails of four separate

ten minute data collections were done. Two of the four were done with classic AdaBoost whereas the other two were done with Modest AdaBoost. Further, under each of these two classifiers, one data collection measured how many false positives were detected, whereas the second data collection counted how many rocking actions went undetected. During all these data collection the researcher counted the number of false positive or false negatives using a handheld thumb counter. This experiment was conducted purely to test the generalization capability of the learnt classifier.

4.6 Experiments

Experiments were carried out for comparing the performance of the classic AdaBoost framework with Modest AdaBoost for the specific tasks of determining

- a. The length of a temporal packet of data needed to effectively distinguish rocking from non-rocking.
- b. The accuracy with which the two classifiers can distinguish between rocking from non-rocking.
- c. The generalization capabilities of the two classifier systems.

To this end the rocking samples collected in Routine A (discussed under Section 4.5.1.1) and Routine B (discussed under 4.5.1.2) were used as labeled positive (rocking) and negative (non-rocking) data for training the AdaBoost classifiers. Data collected under Routine C (discussed in Section??) were used for testing the learnt classifiers. The results from this analysis were used for determining a. and b. above. We varied the packet length on the data stream and determined the recognition rate on the test data. While the packet length was varied, a constant overlap was maintained between successive packets. This overlap was determined empirically to be 0.5 seconds or 50 samples (100 Hz sampling rate). With the ground truth already provided for the test set, we were able to determine the accuracy of the two classifiers.

To determine c ., we resorted to using the data collected in Section 4.5.2. The primary author of the paper used the device to collect false positive and false negative data in order to determine how well the classifiers generalized on the training data. Further, we analyzed the working of the two classifiers in a piece wise manner by breaking down the features into individual sets (Sets 1 through 5 as identified in Table 4.1 and Table 4.2 and Set 6 that included all 14 features) and understanding the functional ability of the classifiers under individual feature sets. This allowed for an in-depth analysis of the workings of the two classifiers. In Section VII, we discuss the generalization capability of the two classifiers by heuristic analysis of the piecewise operational modes.

All our experiments were carried out with the aid of the AdaBoost Matlab library developed by Graphics and Media Lab at the Dept. of Computer Science at Moscow State University¹.

4.7 Results

Figures 4.4 and Figure 4.5 shows the box plot [205] of packet length (T secs) versus recognition rate for classic AdaBoost and Modest AdaBoost frameworks, respectively. The abscissa represents the length of the data stream (in seconds) used for the analysis, while the ordinate represents the recognition rate. Training and testing were all carried out on the data collected as depicted in Section 4.5. The horizontal line inside the box represents the median (second quartile) of recognition rates over the ten subject's data. The lower end of box presents the first quartile (25 percentile) and the upper end of the box represents the third quartile (75 percentile). Thus the box surrounds the center 50 percentile ranges of recognition results. This box is also called the Inter-Quartile Range (IQR = third quartile - first quartile). The dotted extremity represents the minimum and maximum recognition rate under a certain packet length among the ten subjects. Any outlier (an outlier is greater than 1.5 IQR from the median in any direction) is marked by an asterisk.

Table 4.3 presents the results from the experiment carried out to determine the

¹A. Vezhnevets and V. Vezhnevets, GML AdaBoost Matlab Toolbox 0.3, Graphics and Media Lab, Moscow State University, 2007.

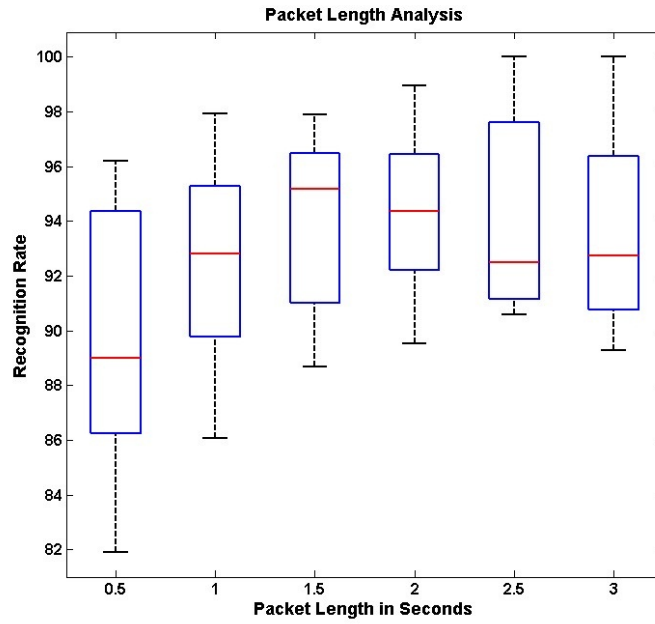


Figure 4.4: Packet length to recognition rate comparison under the classic AdaBoost framework.

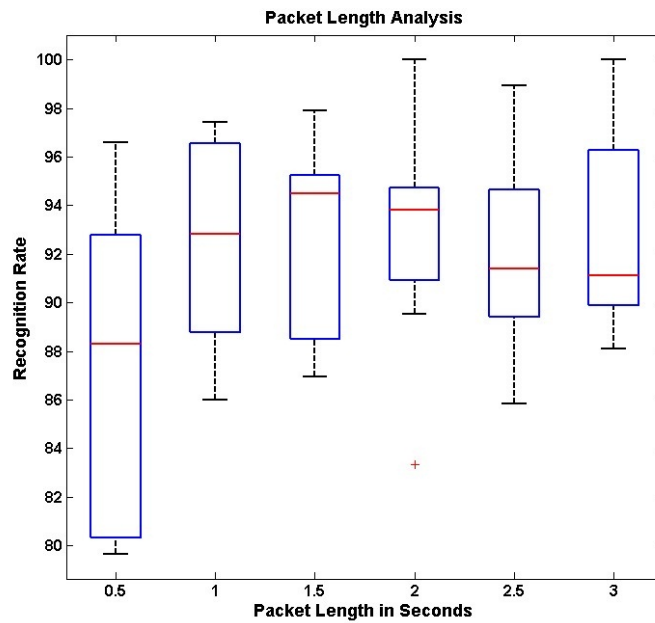


Figure 4.5: Packet length to recognition rate comparison under the Modest AdaBoost framework.

generalization capabilities of the two classifiers. The entries in the table are counts as measured by the researchers of the number of false positives and false negatives counted manually while using the device for body rock detection and feedback. Five trails were carried out of 10 minutes each for determining these numbers. False positives represent the number of times the device falsely gave feedback when the user was not involved in rocking. It is important that this rate be minimal as too many false feedbacks would be discouraging for the user to continue using the assistive aid. The false negative represents the number of times the device did not detect that the user was rocking. This metric could be correlated to the failure of the device to perform its functional task.

Table 4.3: Experiments with naturalistic data

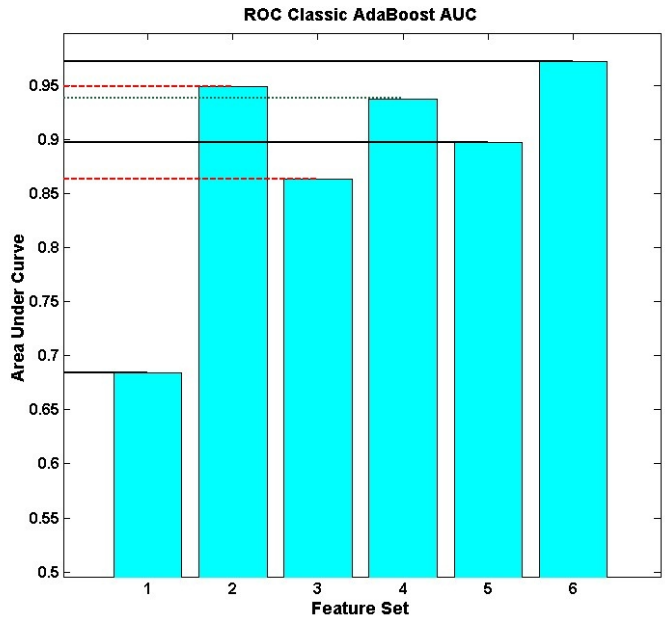
Generalization Capabilities	Classic AdaBoost	Modest AdaBoost
False Positives per Minute ¹	86	44
False Negatives per Minute ¹	20	9

¹ Averaged over 10 minutes

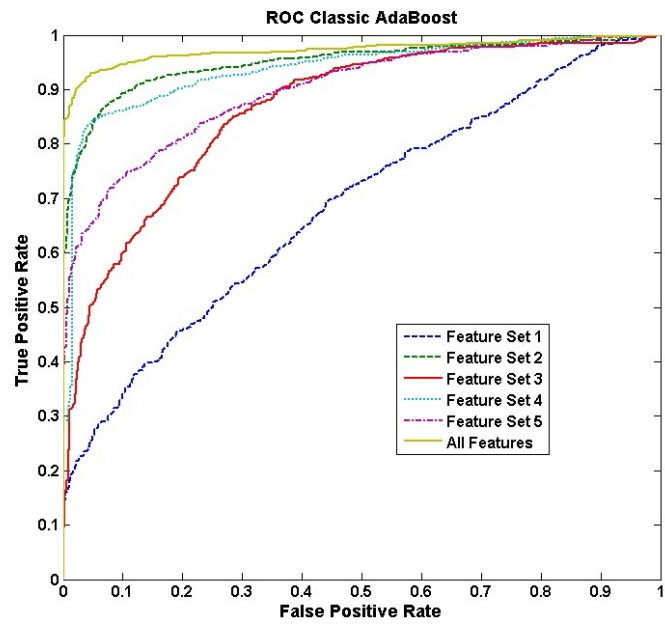
Figure 4.6 and Figure 4.7 shows the piecewise analysis of the classic AdaBoost and Modest AdaBoost frameworks. Subfigure (a) shows the performance of each feature set considered one at a time in detecting body rocking; feature set 6 corresponds to the use of all 14 features together. For example, column 1 in Figure 4.6(a) represents the recognition performance using only temporal mean along x, y and z axis tested on all ten subjects. The bar graph in (a) shows the mean performance rate while the superimposed box plot shows the performance at first, second and third quartile as discussed earlier.

Subfigure (b) represents the Receiver Operating Characteristics (ROC) [206] for the same six feature sets as in subfigure (a). ROC is plotted a false positive rate (FPR) versus true positive rate (TPR). The better the performance, the curve moves towards the (1,1) co-ordinate. For example, in Figure 4.6(b) Set 6 with all features is performing better than feature set 1 as Set 6 curve is closer towards (1,1) while the feature set 1 curve is almost along the diagonal of the plot. The diagonal of the ROC plot represents a recognition rate of 50% i.e. random pick.

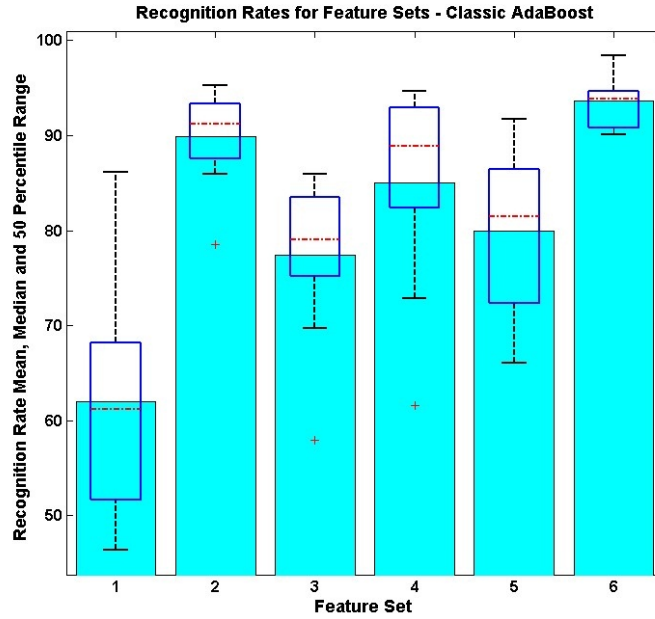
Subfigure (c) is a derivate of the ROC curves in subfigure (b). Each bar in the graph



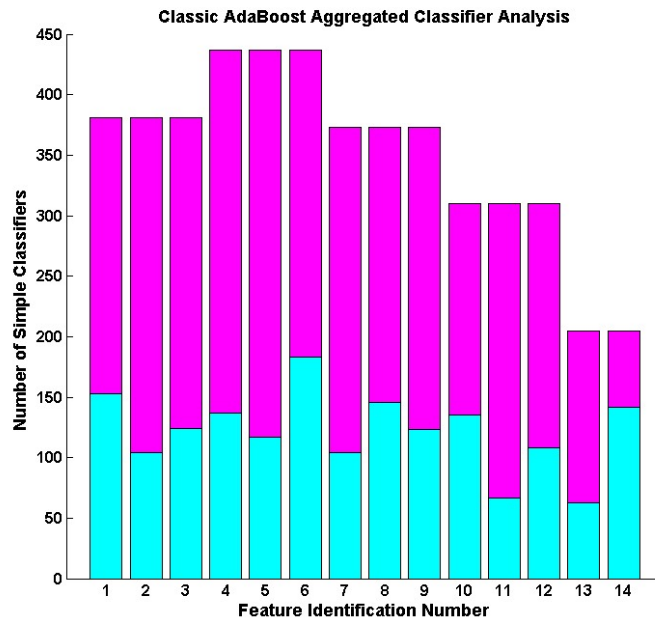
(a)



(b)

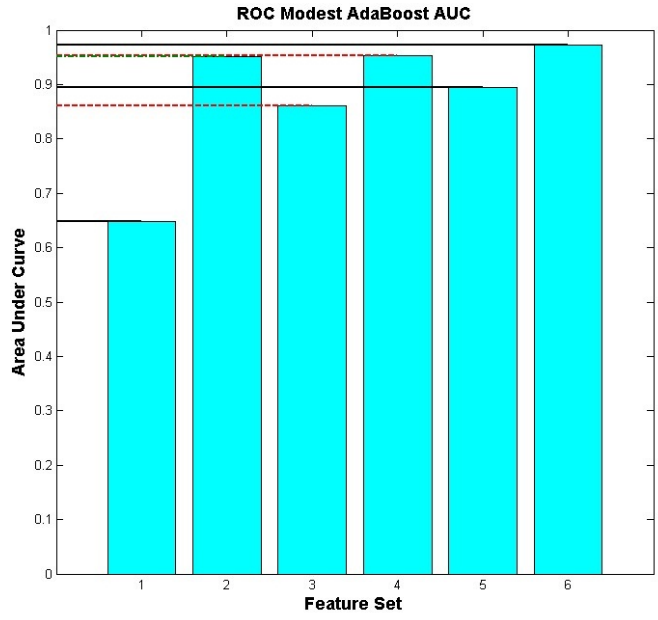


(c)

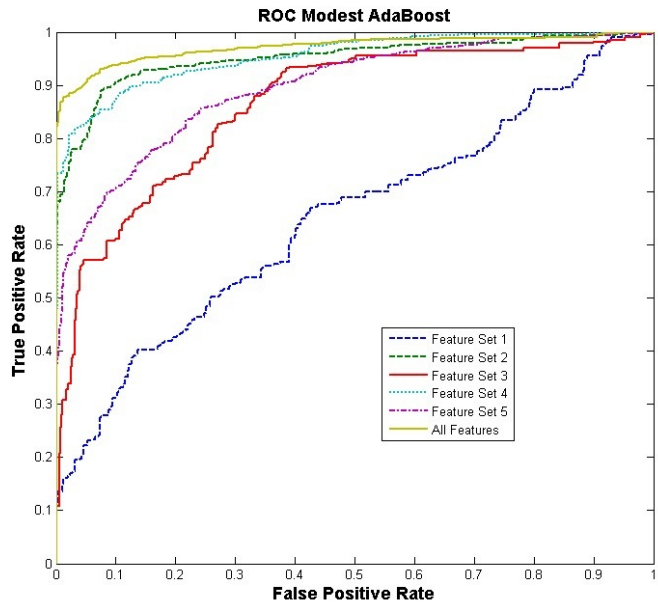


(d)

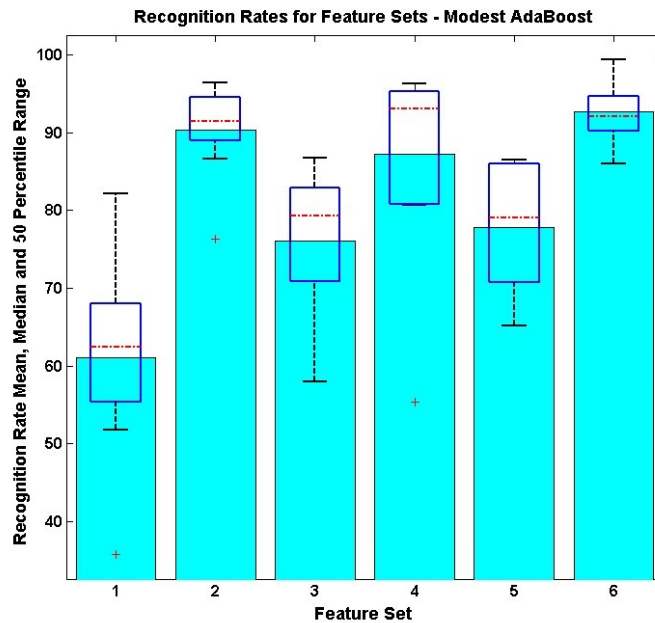
Figure 4.6: Piecewise performance analysis of the classic AdaBoost classifier framework; (a) Recognition rates under use of individual feature sets; (b) The Receiver Operating Characteristics (ROC) under the use of individual feature sets; (c) Area under the curve (AUC) for each feature set as estimated from the ROC; (d) The number of simple classifiers used by the aggregated AdaBoost classifier. Each set and each feature representation in the classifier pool are separately marked. In all the graphs Set 1 through 5 are as explained by Tables 4.1 and 4.2. Set 6 represents a set containing all 14 features from Tables 4.1 and 4.2.



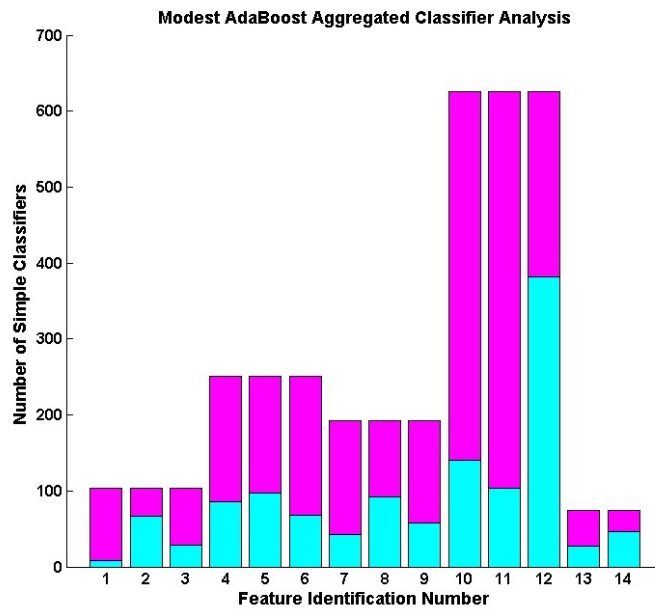
(a)



(b)



(c)



(d)

Figure 4.7: Piecewise performance analysis of the Modest AdaBoost framework; (a) Recognition rates under use of individual feature sets; (b) The Receiver Operating Characteristics (ROC) under the use of individual feature sets; (c) Area under the curve (AUC) for each feature set as estimated from the ROC; (d) The number of simple classifiers used by the aggregated AdaBoost classifier; Each set and each feature representation in the classifier pool are separately marked. In all the graphs Set 1 through 5 are as explained by Tables 4.1 and 4.2. Set 6 represents a set containing all 14 features from Tables 4.1 and 4.2.

is representing the area under the corresponding curve (AUC) in (b). An area of 1 represents an ideal classifier with no false positive or false negatives, while an area of 0.5 represents randomness in the classifier output. AUC can be used to immediately determine the curve with the best performance.

Subfigure (d) is an understanding of how the aggregated AdaBoost classifier is built. As discussed above, AdaBoost classifier uses a collection of simple classifiers to achieve the final classifier. We plotted the number of times a particular feature is being used by the aggregate classifier. Further, the features are grouped into 5 sets corresponding to the five feature sets identified in Table 1. Columns belonging to the same set have the same top count which corresponds to the total simple classifiers used from that set. Each column within the set represents how many classifiers are used on each feature within that set. The count on the individual feature is represented by the bottom color along each column. For example, consider set 4 in Figure 4.7(d), features with identification number 10, 11 and 12 form this set (corresponding to the first order differential power from x, y and z axis of the accelerometer data) and have a top count of 646 simple classifiers. Within the group, the z axis differential power dominates the other two by having a count of 374.

4.8 Discussion of Results

Before discussing the results of the experiments conducted on the accelerometer data, we step back to the first research question that we identified in Section 4.2, Is there any evidence of individuals responding to rehabilitation for reducing stereotypic body rocking behavior? From the Psychology background work presented in Section 4.3, we believe that there is enough evidence that individuals with sensory or cognitive impairment respond to rehabilitation through assistive devices. Specifically, the experiments highlighted in Section 4.3.2.1 support the claim that body rocking can be decreased by providing immediate feedback to the individual.

Regarding the second research question, what is the state-of-the-art technology available to detect and notify individuals of their rocking behavior? We identified the state-

of-the-art motion sensor that is small enough in form factor to become part of one's everyday clothing. Further, we designed this device to be discrete so that the user does not feel any intrusion into their everyday activities. The software can be run on any mobile processing device that the user already carries like a cell phone or PDA. This allows the users to use the device without carrying any additional load. Our solution to this research question caters to the Acceptance design dimension that we identified in Section 4.3.2.1 1.

Focusing on the third research question, Is it possible to build a device that detects body rocking condition and how well can it distinguish body rocking from other functional activities of daily living? We turn our attention to the various results presented in Section 4.7 to prove the efficiency of our proposed method in detecting body rocking and distinguishing it from other non-rocking behavior.

4.8.1 Packet Length, and Detection Efficiency

From Figure 4.4 and Figure 4.7, it is evident that the recognition rates for the two classifiers are comparable and the median recognition rate ranges from 89% to 95%. Based on these numbers, the best performance was achieved at a sample length of 1.5 seconds or 150 samples per packet. Packet length of 150 samples has the highest recognition rate on both the classifiers. Comparing this packet with the 2 seconds packet length or 200 samples per packet, we notice that the 2 seconds packet is very close behind and it has a smaller 1.5 IQR box. Thus, the variance in the recognition rates between 10 subjects is lesser in the 200 samples packet length, implying that the results are more consistent. Further, we noticed that the average natural rocking motion of the 10 subjects was around 27 rocks a minute (i.e. 27 rocks in 60 seconds or 2.22 seconds per rock; this is supported by results from [207]), which implies that a latency of 2 seconds was the closest to the time duration of a single rocking action. As mentioned earlier, all experiments were carried out with an overlap 0.5 seconds or 50 samples between successive packets. Combining these two results, we have

1. *Optimum Packet Length* 2 seconds or 200 samples with 0.5 seconds or 50 samples

overlap between packets.

2. *Best Detection Rate @ 2 seconds packet length $\approx 94\%$ under both classifiers*

4.8.2 *Generalization Capabilities*

From Figure 4.4 and Figure 4.7, it is very difficult to distinguish any performance benefits between classic AdaBoost and Modest AdaBoost. But analyzing Table 4.3, we can notice a dramatic difference in the performance of the Modest AdaBoost when compared to classic AdaBoost. The number of false positives is down from 86 to 44 over a ten minute period. That is, the user receives nearly half less number of false feedback with Modest AdaBoost framework when compared to the classic AdaBoost. This was not evident in the detection tests that were carried out with data collected from Routine C (Section IV - A 3.). We asked the question of why there is an increased performance in Modest AdaBoost and why there is a discrepancy between the test results from Routine C and the naturalistic data capture (Section 4.5.1.3). The answer to these questions lies in the generalization capabilities of the two classifiers. We noticed that most of the false feedback provided by classic AdaBoost occurred while the user was sitting and not rocking. In hind sight, we realized a slight discrepancy in our non-rocking (negative class) data collection. While capturing data under Routine B (as explained in Section 4.5.1.2) the participants were asked to perform various tasks that did not involve rocking to use as negative training set. We realized that most of the participants performed tasks that involved some form of walking or standing activities while they did no activity that involved sitting and not rocking. Thus, just sitting activity was a non-rocking event that was not represented in the training data set. We hypothesize that classic AdaBoost over trained on the non-rocking data while Modest AdaBoost, which is penalized for learning the training set very well, had a better generalization. Extending this heuristic analysis to a more formal analysis, we look at the piecewise performance of the two classifiers. Comparing the ROC curves from Figure 4.6 (b) with Figure 4.7 (b), it can be seen that feature set 2 - Variance and feature set 4 - First Order Differential Power performed the best following Set 2 - All features set. Now comparing Figure 4.6 (d) with Figure 4.7 (d) it can be seen that Modest AdaBoost distributed it simple classifiers such that

there were more classifiers representing the two feature sets 2 and 4. On the other hand, the classic AdaBoost's distribution of simple classifiers is unexplainable as feature set 1 - Mean - seems to have received more representation than set 4. Mean had the worst performance as an individual feature set as can be verified by the ROC curve that comes closest to the diagonal on the plot hinting that the performance is barely above random guess. Contrasting this with Modest AdaBoost selection, Mean is in the bottom two sets among the five feature sets. This bad performance of Mean as a feature set can be understood by looking at the graph shown in the first row and last column of Table 1. It can be seen that the Mean acceleration values between rocking and non-rocking are not significantly different. Table 4.1 Row 2 and Table 4.2 Row 1 highlights the capabilities of Variance and First Order Differential Power in distinguishing rocking from non-rocking. This is further confirmed by the ROC graph.

Feature 4 having the highest distribution of simple classifiers under Modest AdaBoost (Figure 4.7 (d)), within this feature set we can see that the highest number of simple classifier is assigned to feature 12 which corresponds to First Order Differential Power on z axis. As can be verified from Figure 4.3, the best distinguishing character between non-rocking and rocking patterns seems to be the transformation of a random signal pattern on z-axis to a deterministic sinusoidal waveform. If this can be the true identity of the rocking data stream, feature 12 would capture it in the best possible manner by measuring the power in the first order differential of the temporal signal. Using this feature as the most reliant feature would provide a good basis to support the final classifier selected by Modest AdaBoost.

We are now ready to answer the last research question stating that the use of approximately 2 seconds (or 200 samples @ 100 Hz sampling rate) packet length used with a learning framework biased towards generalized learning (like Modest AdaBoost) would be a good assistive technology solution for detecting and giving feedback towards stereotypic body rocking. We can extend the same argument to other body mannerisms that involve any form of repetitive body part movement.

In this chapter, we have addressed the topic of detecting stereotypic body mannerisms, specifically body rocking, and propose a technology solution for providing an assistive technology that may reduce or control body rocking. We have discussed the hardware and software components of the proposed system in detail and offer a thorough analysis on the learning framework that provides generalization benefits to allow this framework to be extended to detection of any body mannerism. Investigations are in progress to determine how incoming samples of acceleration data can be labeled automatically by the system based on the AdaBoost classifier's classification confidence metrics. This would provide opportunity for self-learning [208] modes where the device can readily understand and learn data points that were not available in the training set. Combining such self-learning into a generalized learner would provide immense opportunities for not only body mannerism detection, but for solving future data mining problems where typical lab setting training data collection would just not be sufficient to train a robust classifier.

SENSING FACIAL EXPRESSIONS IN DYADIC INTERACTIONS

As described in Chapter 3, the evidence towards facial expression recognition precedes other non-verbal cue sensing and delivery. This is supported by the argument that most part of the non-verbal cues occur through visual facial mannerisms as described in Figure 1.1. The face encodes a lot of information that is both communicative and expressive in nature. Unfortunately, the face is a very complex data generator and the encodings on the face are very individualistic in nature. Evolving computing technologies have been focused on developing solutions towards understanding the nature of facial mannerisms and gestures, but most of this research has been focused on the development of sensors and algorithms that understand user's emotional state for a human-machine interaction scenario. Such interactions are mostly unilateral in nature and focused primarily on developing technologies that will allow the machine to interpret the user's emotional state. That is, the machine becomes the primary *consumer* of the affective cues. But from the perspective of an assistive technology the extracted facial expression primitives have to be augmentations that enrich human-human interpersonal interaction, where the machines not only interpret communicator's affective state, but also deliver affect information through novel affect actuators to the user of the technology. That is, the technology needed in a social interaction assistant is not a consumer, but a *mediator*.

Affect information is causal in nature and understanding what the expression or mannerism means requires an understanding of the context in which it took place. State-of-the-art computational models developed towards understanding context are very simplistic and performs only nominally even under very well controlled laboratory conditions. Contrary to such a setting, assistive technologies provide some respite to the complexities by allowing the user to make all the cognitive decisions. That is, while human computer interfaces need to mimic sensing, cognition and delivery, assistive technologies for people who are blind have to look at sensing and delivery alone and include the human cognition in the loop. This requires precise sensing of the facial and head movements while delivering as

much information back to the user as possible through technologies that do not overload the user with information but provides just the right level of information to allow them to cognitively process this information.

Facial expressions are of utmost importance when humans interact on a one on one basis as explained in Section 1.2.1.1 of Chapter 1. That is, facial expressions become very relevant when the person who is blind is involved in a dyadic interaction with his/her sighted counterpart. To this end, it is essential that the technologies developed towards enabling the user to access facial expressions be a) able to precisely sense facial movements, b) able to deliver facial expression information seamlessly to the user through non-visual modality. Below we consider the challenges of dealing with facial expressions in a dyadic interaction scenario.

5.1 Design Considerations

The human face is very dynamic when it comes to generating important non-verbal communicative cues. Subtle movements in the facial features can convey great amounts of information. For example, slight opening of the eyelids conveys confusion or interest, whereas a slight closing of the eye lids conveys anger or doubt. Thus, the human face can be considered to be a very high bandwidth information stream, where careful design considerations need to be taken into account if this data has to be encoded optimally and effectively through other modalities.

In the past, most researchers and technologists have resorted to auditory cueing when information has to be delivered to persons with visual disabilities; but there is a strong growing discomfort in the target population when it comes to overloading their hearing. People with visual disabilities have a natural tendency to accommodate for the lack of a primary sensory channel by relying on hearing. For example, with the aid of ambient noise in a room, they can gauge approximately how big a room is. Thus, when designing assistive devices aimed at social aid, we need to carefully consider how to deliver high bandwidth data streams to users relating to the facial movements of interaction partners. Touch or

haptic based delivery is a growing area of research which is relatively underutilized, except for Braille. In the next chapter, we explore the use of vibrotactile cueing (on the human hand through a Haptic Glove) towards delivering facial expression data to the user. Below we provide the overall application scenario for the proposed dyadic interaction assistant the details of which are presented in this and the next chapter.

5.2 Dyadic Interaction Assistant - Proposed Solution

Dyadic interactions represent a large portion of social interactions between individuals; and during dyadic interactions, it is very important to assess the communicator's face, head and body-based gestures and mannerisms. The dyadic interaction assistant is meant to convey important facial and head mannerisms of the interaction partner that might correlate to communicative gestures like head nod, head shake, doubt, anger, frustration, happiness, interest, etc. See Figure 5.1 for a typical dyadic interaction setting. The device incorporates an automated table top face tracker which acts as the input to the system while allowing any extracted data to be delivered on a wearable haptic display called the Haptic Glove.

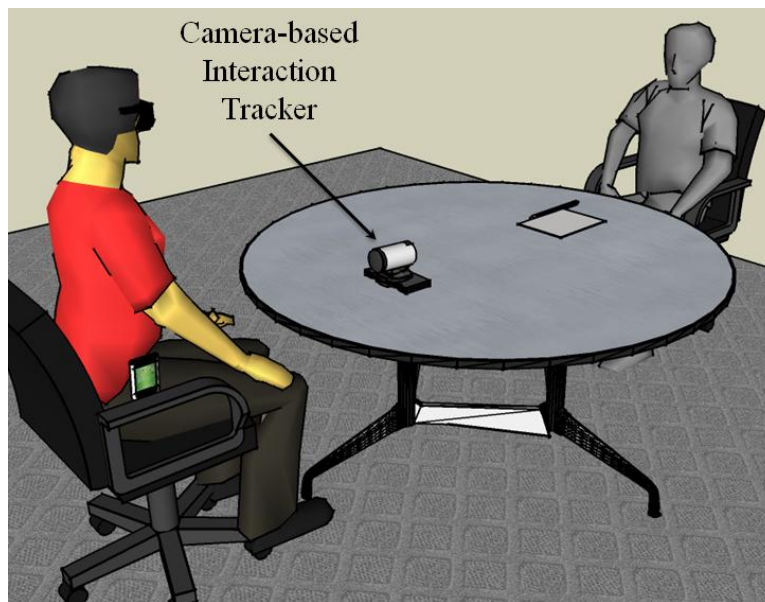


Figure 5.1: Typical use of the dyadic interaction assistance scenario, a third person perspective on the use case scenario.

The dyadic interaction assistant was designed to be a compact device that can be

carried into meetings where an individual who is blind or visually impaired could place it in front of his/her interaction partner. The device, as shown in Figure 5.2, consists of a micro pan-tilt mechanism that is controlled from a PDA like computing platform. Real-time face detection (as explained in the previous section) tracks the face of the interaction partner and captures only the face image for further processing. We use the FaceAPI, a commercially available facial feature tracking package to determine the locations of all the facial features including eyelids, eye brows and the mouth. Figure 5.3 shows a typical output of the FaceAPI software.

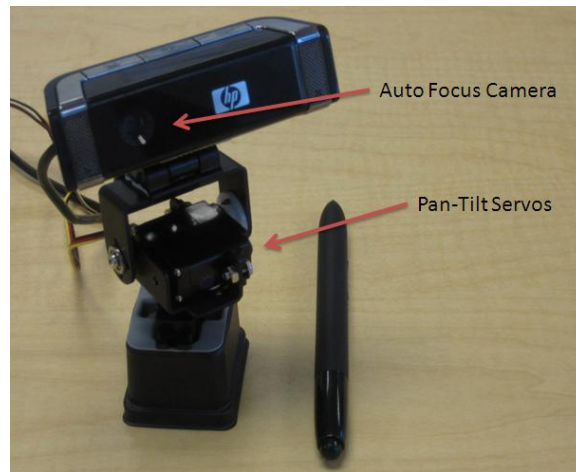


Figure 5.2: Face tracker with an auto-focus camera and a micro pan-tilt mechanism.

This chapter focuses on the sensing of facial expressions using the proposed dyadic interaction assistant. Since facial expressions are all visual in nature, the proposed solution uses computer vision solutions to track the face of the interaction partner and extracts all the necessary facial feature information for decision making and delivery. The following section gives a summary of the current state-of-the-art in computer vision based facial expressions recognition algorithms. Following this discussion, the proposed methodology for extracting the expression primitives are detailed.

5.3 Facial Expression Recognition - State of the art

Computer based facial expression recognition is mostly centered around the process of identifying the various facial mannerisms and gestures of individuals through the use of mostly



Figure 5.3: Typical use of the dyadic interaction assistance scenario, a third person perspective on the use case scenario.

electrooptical cameras. While recent developments relate to their recognition using infrared [209], thermal [210] and face electromyogram (EMG) [211], in this chapter, we restrict to the discussions of EO camera based facial expression recognition as the primary sensor on the dyadic interaction assistant is restricted to the use of an auto-focus Complementary Metal Oxide Semiconductor (CMOS) camera.

Most facial expression recognition algorithms treat the underlying problem as a classification problem and group spatio-temporal facial actions into predetermined bins of specific facial gestures or mannerisms. Popular classifications include the Ekman [212] grouping of 6 basic facial expression of emotion, including Happy, Sad, Surprise, Fear, Anger and Disgust. These classes are distinguished from the otherwise Neutral expression of the individual's face.

Table 5.3 below shows a compilation of the state-of-the-art in facial expression recognition algorithms. Note that the table is an extension of the 2009 work that was published by Zaho et al. [213].

Table 5.1: State-of-the-art facial expression recognition algorithms and their performance.

Exp: Spontaneous(S) / Posed expression (P); Per: Person Dependent(P) / Independent (I); Class: Number of classes; Sub: Number of subjects; Type: Data Type: Video (V) / Image(I); % Acc: Percentage Accuracy; ?: missing entry.

Ref.	Features	Classifier	Performance					
			Exp	Per	Class	Sub	Type	% Acc
[214]	AAM	SVM	S	I	2	21	?	81
[215]	Gabor	SVM + HMM	S	I	3 AUs	17	V	98
[216] [217]	Gabor	AdaBoot SVM	S P	I	17 AUs	119 + 12	I	93 + 90.5
[218]	12 motion units	Tree DBN HMM	P	D I	6	5 + 53	V	66.5 + 73.2
[219]	Shape Models, Gabor	LDC	S	I	3 AUs	21	I	76
[220]	24 facial points	DBN	P	D	6	30	V	77
[221]	Intensity	NN	P	?	7	?	I	68
[156]	Shape fea, Optic flow	C4.5 Bayes Net	P	?	8	4	I	100
[222]	FAPs	Neurofuzzy network	S	I	3	?	I	78
[223]	Shape fea	DBN	S	?	2	8	V	95.3
[224]	Facial and head gesture	GP SVM HMM NN	S	?	2	8	V	86
[225]	Pixel diff of mouth	GP SVM HMM NN	S	I	2	24	V	79
[226]	Intensity of face	Decomposable model	P	I	6	8 + 16	V	61
[227]	Gabor	AdaBoost SVM	S	I	2	26	V	72
[228]	AAM	SVM	S P	I	AUs	100	?	95
[229]	Facial profile	Rule-based	P	I	27 AUs	19	V	86.3
[230]	Frontal profile facial points	Rule and case based	P	I	9	8	I	83
[231]	12 motion units	kNN	S	I	4	53 + 28	V	93 + 95
[232]	Gabor	AdaBoost DBN	P	I	14 AUs	100 + 10	I	93 + 93
[233]	Motion history	SNoW kNN	P	I	15 AUs	19 + 100	V	61 + 68
[234]	8 facial points	Gentle Boost SVM	S P	I	2	27 + 32 + 65	V	90
[233]	20 facial points	Gentle SVM	S P	I	2	52	V	94
[235]	Shape fea, Intensity	NN	S	?	7	14	I	84
[236]	3D surface	LDA	P	I	6	60	I	83
[237]	Geometric ratio	GMM	P	I	4	47	I	75
[238]	Harr	AdaBoost	P	I	11 AUs	?	I	92
[216]	Intensity	kNN HMM	S P		6	97 + 21	V	90.7 + 82
[239]	Texture with LPP	SVDD	S	D	2	2	I	87

5.4 Design Considerations for Dyadic Interaction Assistant

As seen from the Table 5.3, facial expression recognition algorithm is a mature area with various solutions being suggested towards myriad applications ranging from biometrics to human computer interfaces. In contrast, the requirements in the proposed application is focused towards mediating human-human communication and thus imposes important constraints like,

- *Real-time interpretation of the facial expressions and gestures:* People who are blind, or disabled in general, prefer not to rely on technologies that do not respond to their needs immediately and reliably. Survey of early adoption of assistive technologies show that the technology need not accomplish a large range of tasks, but the tasks being performed should be executed reliably and in useful time [240]. Range of facial expressions range from a few hundredths of a second to couple of seconds at max ([5], Page 322). It is essential that the facial expression recognition algorithms be able to respond within this time frame to enable seamless interactions with sighted counterparts.
- *Ability for the user to choose varying levels of facial information of their interaction partner:* In our interactions with the user community, it has become evident that each user has his/her own requirement of the details of the facial information that they receive from their interaction counterparts. For example, in a professional setting, individuals preferred to receive information that was delicate and down to detail, while in a personal setting, they preferred that the device does not overwhelm the user with facial interaction information.
- *Ability to switch between subject-dependent and subject-independent facial expression recognition:* Social interactions vary between personal interactions with friends and family to professional meetings with strangers that the user may meet only once. An important aspect of social interaction is the ability to relate to an individual's personal mannerisms and gestures. This is especially true in personal interactions and

behavioral psychology studies shows the inherent need for individuals to recognize and reenact their interaction partner’s gestures (popularly termed as the *Chameleon Effect*). From the sensing perspective, thus it is important that the algorithm be able to adapt to generic facial expression recognition, while also be able to adapt to subtle changes in a single subject’s facial movements. That is, the algorithm requires the benefit of switching between subject-dependent and subject-independent facial expression and gesture recognition.

In the following section, we describe a recently proposed framework for facial expression recognition, termed as the Temporal Exemplar-based Bayesian Network, which demonstrates the ability to satisfy the needs itemized above. We first describe the framework and provide justification for its suitability to address the identified requirements.

5.5 Temporal Exemplar-based Bayesian Network (TEBN) for Facial Expression & Gesture Sensing

A Bayesian Network is a interdependency graph representing the influence of various events on a desired outcome. A Bayesian Network is the extension of the Bayes Rule of conditional probabilities to multiple acyclically interdependent variables. For example, consider the problem of estimating the probability of a lawn being wet, under the possibilities that it could have rained, or the sprinklers could have turned on. The probability of wet grass $P(w_g)$ then depends on the probability of raining $P(r)$ and the probability of sprinkler turning on $P(s)$. Thus, the joint probability $P(w_g, s, r)$ can be written as

$$P(w_g, s, r) = P(w_g/r, s).P(s).P(r) \tag{5.1}$$

Further, most sprinklers can be assumed to be dependent on whether it rained or not, there by making the probability $P(s)$ dependent on the rain. Incorporating this into the above equation,

$$P(w_g, s, r) = P(w_g/r, s).P(s/r).P(r) \quad (5.2)$$

Graphically, the above equation can be represented as shown in Figure 5.4,

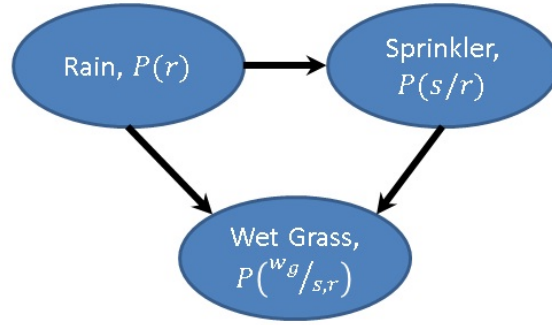


Figure 5.4: Example Bayesian Network.

Typically, the various conditional probabilities associated with the graph are determined through statistical analysis of observations. Most often, the analysis results in inductive models that generalize the training observations to all test cases. In contrast, the Exemplar-based Bayesian Network determines the conditional probabilities transductively on a case-by-case basis of the test set. Thus, the network is dynamic, allowing the conditional densities to transform based on the input test condition. Also, the network accounts for the temporal transmission of decision probabilities from time instance t to $t + 1$. This transmission is also modeled into the network through a Markovian process. The details of the implementation follows, as described in [241].

Given any facial expression image sequence $I(t)$ with M number of frames, the goal of TEBN is to provide a posterior conditional probability of the sequence belonging to a possible set of N facial expression labels Y_i . N depend on the application under consideration. Thus, the conditional probability can be defined as

$$P(Y_i/(I(t), H(t))) = \frac{P(I(t)/(Y_i, H(t)))P(Y_i/H(t))}{P(I(t)/H(t))} \quad (5.3)$$

If we assume that the image sequence $I(t)$ is independent of the past history $H(t)$, given that we know the expression label Y_i and introducing the exemplar layer $L_i(t)$ as all the knowledge from which the labels Y_i will be derived, the likelihood in the equation 5.3 can be rewritten as,

$$P(I(t)/(Y_i, H(t))) = P(I(t)/L_i(t))P(L_i(t)/Y_i) \quad (5.4)$$

Thus,

$$P(I(t)/(Y_i, H(t))) = \frac{P(I(t)/L_i(t))P(L_i(t)/Y_i)P(Y_i/H(t))}{P(I(t)/H(t))} \quad (5.5)$$

Neglecting the scaling factor in the denominator, the network can be defined by three layers, namely, *Observation Layer*, *Exemplar Layer*, and *Prior Knowledge Layer*, as shown in the Figure 5.5. The three layers shown in the figure contribute towards the final decision on the nature of the facial expression. Note that this example represents a classification problem where the six basic human expression (Happy, Sad, Surprise, Anger, Fear and Disgust) are considered to be the prime focus. In the later sections, we describe how this framework can be adopted to suit dyadic interaction assistance requirements.

5.5.1 The Observation Layer

As described earlier, we use the results from the FaceAPI software as the inputs into the facial expression recognition algorithm. As shown in Figure 5.6, 36 points returned by the software are used as the input observation vector, $X(t)$ into the TEBN. The length of this observation vector, N , is $2 \times 36 = 72$, corresponding to (x, y) coordinates of the 36 facial fiducial. Structurally, the 72 data points are arranged as a linear vector $X(t)$ that represents the structural configuration of the facial features at time instant t . As time progresses, the vector X evolves into an expression or gesture. $X(t)$ extracted at each frame is used for choosing K examples from each of the N expression classes, 6 in the example described in Figure 5.5.

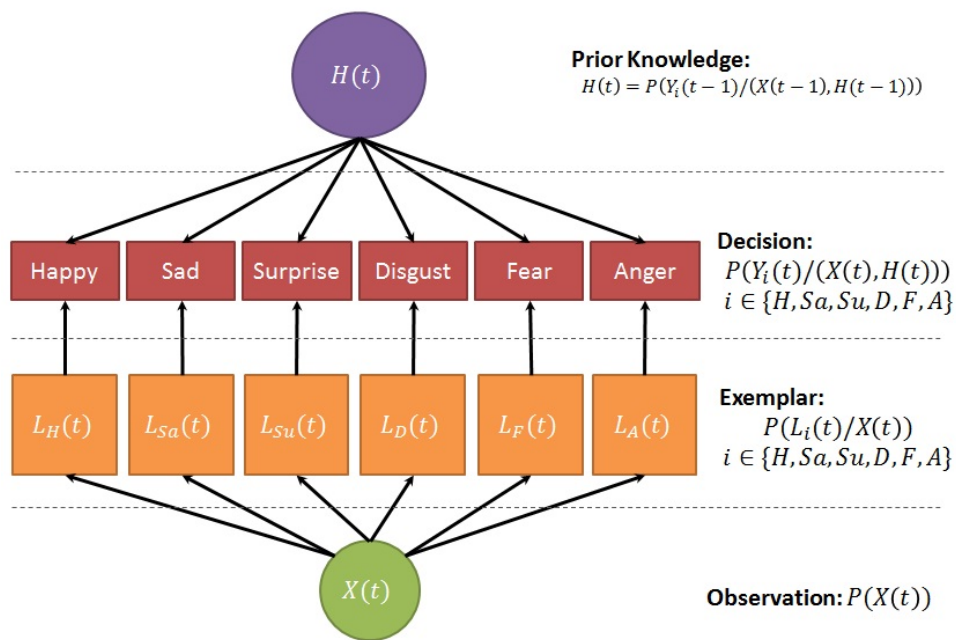


Figure 5.5: Temporal Exemplar-based Bayesian Network for facial expression recognition.

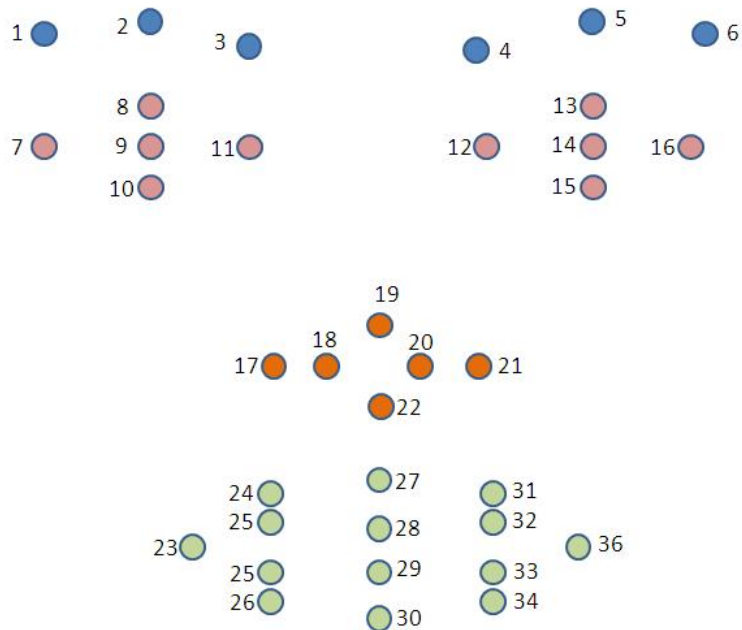


Figure 5.6: 36 Facial fiducial points tracked with FaceAPI software. Both x and y coordinates from all 36 points are used for facial expression recognition.

Given that all training data is available at sequence of images, $I(t)$, we choose to represent each frame of the image just by their FaceAPI tracked points. Thus, in the above equations, $I(t)$ can be conveniently replaced with $X(t)$, the 72 point vector of tracked points. Note that this is done only to reduce the computational load. Later, in the experiments section, we demonstrate the possibility of representing $I(t)$ through other image features.

Thus, equation 5.4 can be rewritten in terms of the extracted features $X(t)$ as,

$$P(X(t)/(Y_i, H(t))) \propto P(X(t)/L_i(t)) P(L_i(t)/Y_i) P(Y_i/H(t)) \quad (5.6)$$

5.5.2 The Exemplar Layer

The Exemplar Layer represents the aggregate knowledge for every test facial expression to be evaluated by the TEBN. As the name suggests, given a test sample of image sequence, $I_t(t)$, for which we need to estimate the facial expression label, Y_i , the exemplar layer is constructed dynamically to best represent the test data. Thus, the exemplar layer $L_i(t)$ is chosen from the training pool, $E(t)$, by comparing $I_t(t)$ with every candidate image sequence in $E(t)$. For each of the expression labels, $Y_i, i = \{1, \dots, N\}$, k nearest neighbor points are chosen in $E(t)$ to represent the test data. We chose Euclidean distance between the test observations $X_t(t)$ and all observations $X_E(t)$ from the training pool $E(t)$. Once the k nearest neighbors have been identified for the given test sequence, the Bayesian Network is developed by representing the test point $X_t(t)$ as the weighted average of L_{ij} training points chosen from the training pool. The examples $L_{ij}, i = \{1, \dots, N\} \& j = \{1, \dots, k\}$ represents k examples taken from N expression classes. Figure 5.7 shows the steps involved in developing the Exemplar Layer for any given test point $X_t(t)$ and highlights the problem of finding w_{ij} such that

$$X_t(t) = \sum_{i=1}^N \sum_{j=1}^k w_{ij} L_{ij} \quad (5.7)$$

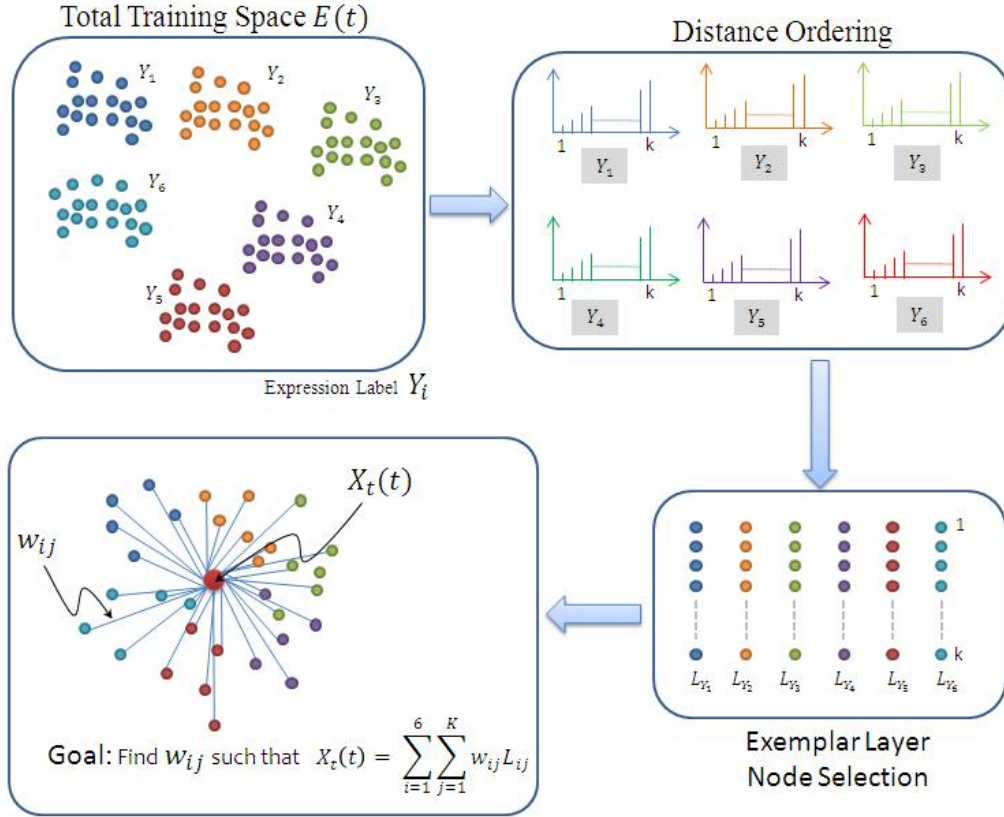


Figure 5.7: Deriving the exemplar layer of the TEBN based on every test point $X_t(t)$.

5.5.2.1 Computing $P(X(t)/L_{ij}(t))$: Representing the test data in terms of existing examples

If any test point $X(t)$ can be represented as the sum average of a set of chosen example training data, in our case $L_{ij}(t)$ (chosen based on k nearest neighbors of $X(t)$), the likelihood $P(X(t)/L_{ij}(t))$ can be computed as

$$P(L_{ij}(t)/X(t)) = \frac{\sum_{j=1}^k w_{ij}(t)}{\sum_{i=1}^N \sum_{j=1}^k w_{ij}(t)} \quad (5.8)$$

and,

$$P(X(t)/L_{ij}(t)) = \frac{P(L_{ij}(t)/X(t))P(X(t))}{P(L_i(t))} \quad (5.9)$$

where, $P(L_i(t))$ can be considered as a proportionality constant and the probability of occurrence of any $X(t)$ is considered to be equal. Thus,

$$P(X(t)/L_{ij}(t)) \propto P(L_{ij}(t)/X(t)) \quad (5.10)$$

Determining $w_{it}(t)$ given test point $X(t)$ and examples $L_{ij}(t)$: In [241], the problem of determining $w_{ij}(t)$ is posed as a entropy maximization problem,

$$\text{Maximize} \quad \left(- \sum_{i=1}^N \sum_{j=1}^k w_{ij}(t) \ln(w_{ij}(t)) \right) \quad (5.11)$$

$$\text{Subjectto} \quad \sum_{i=1}^N \sum_{j=1}^k w_{ij}(t) = 1$$

$$\sum_{i=1}^N \sum_{j=1}^k w_{ij}(t)L_{ij}(t) = X(t)$$

(5.12)

The above problem is a convex optimization problem with the local maxima being the same as the global maxima. Appendix B shows the solution to be problem using Newton's method and the steps shown in subsection B.2.1 of Appendix B shows the algorithm for determining the weights w_{ij} associated with L_{ij} .

5.5.2.2 Computing $P(L_i(t)/Y_i)$: Determining the prior probability of chosen examples w.r.t the complete training data

The Exemplar-based Bayesian model assumes that all the known knowledge of the problem is embedded within the complete training set $E(t)$. Thus, given a subset of the examples, $L_{ij}(t)$ that is extracted from the training set, the likelihood of a given expression label, Y_i can be obtained as the ratio of likelihood of occurrence of the selected training samples over all the training samples in the database. This can be achieved through,

$$P(L_i(t)/Y_i) = \frac{\sum_{x \in L_i} P(x)}{\sum_{x \in E} P(x)} \quad (5.13)$$

The probability $P(x)$ can be determined through Kernel Probability Density Estimation as

$$P(x) = \frac{1}{nh^D} \quad (5.14)$$

Assuming ϕ to be a Gaussian Density function,

$$\phi_h(z) = \frac{1}{(2\pi)^{D/2}} \exp\left\{-\frac{1}{2}z^T z\right\} \quad (5.15)$$

where, h represents the window width and an optimal value for this can be found (See [242] for more information) as,

$$h_{opt} = \sigma \left\{ \frac{4}{n(2D+1)} \right\}^{\frac{1}{D+4}} \quad (5.16)$$

σ^2 is the average marginal variance on the D dimensions of the data in the set E .

5.5.3 The Prior Knowledge Layer

Facial expressions evolve over time and the rate at which one displays a certain facial expression varies from one individual to another. Further, some of the facial expressions show commonalities that exhibit temporal correlations. For example, consider the expressions Fear and Surprise as shown in the Figure 5.8. The movements on the upper half of the face are very similar in the two expressions, but not alike. The differences in the final expression label depends on the curvature and the opening of the mouth. Thus, a temporal evaluation of any facial expression requires that the frame train be analyzed as an evolving probability confidence over time history $H(t)$.



Figure 5.8: Comparison of Fear and Surprise facial expressions.

5.5.3.1 Computation of $P(Y_i/H(t))$: Temporal propagation of expression probabilities

The temporal nature of the TEBN is embedded in the propagation of probabilities across time (video frames). That is, propagation of expression probability from time t to $t + 1$, can be represented as a function of the transition probability given the frame history,

$$P(Y_i/X(1, \dots, t)) = \sum_{j=1}^N P(Y_i/(Y_j, X(1, \dots, t-1)))P(Y_j/X(1, \dots, t-1)) \quad (5.17)$$

Considering a Markovian nature of probability propagation, we can rewrite the above equations as,

$$P(Y_i/X(1, \dots, t)) = \sum_{j=1}^N P(Y_i/Y_j)P(Y_j/X(1, \dots, t-1)) \quad (5.18)$$

In the above equation, the propagation is centered only on the transition probability $P(Y_i/Y_j)$. It is difficult to determine the transition probabilities before hand as they are dependent on the length of the video stream and also dependent on the duration of the expression in question. Hence, the transitions are determined while the expression evolves in the video frames. To this end,

1. Define a transition matrix T of size $N \times N$, to describe transitions between any of the N expressions.
2. Initialize T to 1s.
3. At at given time instance, t , the transition probability, $P(Y_i/Y_j)$ can be obtained as

$$P(Y_i/Y_j) = \frac{T_{ji}(t)}{\sum_{k=1}^N T_{jk}(t)} \quad (5.19)$$

4. On an iterative basis update T matrix based on a frame-by-frame transition from expression i to j . That is, increment the entry T_{ij} by 1, if the $t - 1$ frame is assessed to be expression Y_i , through $P((Y_{X(t)} = Y_i)/X(1, \dots, t - 1))$ and the expression label transitions to Y_j at time frame t , through $P((Y_{X(t)} = Y_j)/X(1, \dots, t - 1))$.

5.6 Discussion of Design Considerations: TEBN perspectives

Having discussed the TEBN framework, we highlight the importance of using TEBN for addressing the three important design constraints that were highlighted in Section 5.4. Below we highlight the importance of TEBN framework in addressing this design considerations.

- *Real-time interpretation of the facial expressions and gestures:* The TEBN framework enables real-time performance. See the experiments section to determine the time taken towards each expression assessment.
- *Ability for the user to choose varying levels of facial information of their interaction partner:* The number of decision levels that the TEBN can have is controlled by the design of the exemplar layer. In the discussions a single layer (6 expression) TEBN was discussed, but the same can be extended to incorporate more classes of expression bases like the Actions Units (AU), various intensity of AUs or the raw movement data itself. That is, the number of elements representing N in the above discussions can be extended based on the classifications desired by the application, scenario and user choice.

- *Ability to switch between subject-dependent and subject-independent facial expression recognition:* Subject dependence and independence can be easily controlled through the determination of the example pool from which the data is being selected for building the TEBN. If the identity of an interaction partner can be established before the expression analysis, it is possible to build $X(t)$, the test vector as a sum of expressions that belong to that specific individual, L_i^r , where $i \in \{1, \dots, N\}$ represents the N expressions and $r \in \{1, \dots, R\}$ represents the total number of individuals in the database. On the other hand, if the interaction partner's identity cannot be established, all expression examples independent of the subject identity can be used for building the exemplar layer for $X(t)$.

5.7 Experiments

5.7.1 Data

All experiments were carried out on the Cohen-Kanade Facial Expression Database [243] version 1. Version 1 (the original or initial release (Kanade, Cohn, & Tian, 2000)) includes 486 sequences from 97 posers. Each sequence begins with a neutral expression and proceeds to a peak expression. The peak expression for each sequence is fully FACS (Ekman, Friesen, & Hager, 2002; Ekman & Friesen, 1979) coded and given an emotion label. The emotion label refers to what expression was requested rather than what may actually have been performed. In order to validate the emotion that was displayed, the final expression images of the 486 sequences were displayed to unbiased participants through a web portal. Each emotion image was displayed to the user and requested to label them as belonging to one of 7 classes, *Happy*, *Sad*, *Surprise*, *Fear*, *Disgust*, *Surprise* and *Neutral*. During the screening process, it was noticed that some of the subjects did not display expressions as requested (either displayed neutral expression through the sequence or displayed exaggerated facial mannerism unlikely of requested emotion). Hence all the participants who labeled the expression images were required to label on a scale of 1 to 5, the genuineness of the expression displayed, 1 being least correlation to genuine expression and 5 being most genuine. Based on the user responses, we identified 182 sequences where subjects displayed

emotions that were named as being genuine up to a scale level of 4 or 5. These sequences were then used for testing the TEBN.

In this chapter, we discussed the use of a Temporal Exemplar-based Bayesian Network to determine the facial expressions of the interaction partner. Given that the dyadic interaction assistant is able to extract these facial expression cues, it is thus required to deliver this information back to the users who are blind. This is a challenge in itself as the very high bandwidth information delivered by a human face could easily overload the user. As explained before, the next chapter investigates the use of the human perception of touch to develop novel haptic interfaces for delivering highly dynamic facial expressions.

DELIVERING DYADIC INTERACTION CUES THROUGH VIBROTACTILE
STIMULATIONS

Over the past few decades, the environments in which humans live and operate have become increasingly information rich. Audio and visual modalities have been occupied by more than one source at the same time. Increasing audio-visual stimulations in the human's surroundings have increased the need for divided attention, which competes with the need for selective and sustained attention to complete a task at hand [244]. While audio and video have evolved as a medium for immersing humans in rich sensory experience, touch [245], taste [246] and smell [247] have only recently been considered for sensory augmentation and substitutions - note the subtlety between augmentation and substitution. When augmenting the already utilized sensory channels, the newer medium is not in demand to reproduce all of the information, but only enrich the already delivered experience. On the other hand, substitutions have to deliver information that was once being provided by a certain sensory channel on a newer medium, while maintaining similar or lesser cognitive load.

Vision is the primary sense organ for most mammals and for a few primates, including humans, trichromatic vision is so highly evolved [248] that a major portion of the neuronal pathway in the brain is dedicated to sensing, perceiving and cognizing visual stimuli. This allows the human vision to process high intensities of data that stimulates the eyes - Koch et al. estimated that human eyes, with 106 ganglion cells, could transmit up to 10 Mbits/s [249] to the brain. Hence substituting vision with any other sensory channel is a challenging task that requires appropriate design taking into account the high intensity of data generated by visual stimuli.

In this chapter, we present a visuo-haptic sensory substitution device that intends to replace visual channel with somatosensory vibrotactile stimulations for delivering the facial expression data that was derived from the dyadic interaction assistant, as described

in the pervious chapter. A prominent visuo-haptic sensory substitution system is the TVSS (Tactile-Visual Sensory Substitution) [250] that substitutes visual data into a 400 point tactile array worn on the back of the user. A similar effort by Rahman et al. [251] focused on delivering facial emotions of interaction partners using vibrators installed on a chair such that the user's back is in direct contact with the vibrators. Both of these technologies have the obvious disadvantage that the user is restricted to a seated position with immobile augmentations. Recently we have explored the use of vibrotactile technologies on a belt like form factor for delivering direction and distance information [252]. Unfortunately, the waist (combined with a belt like form factor) did not prove suitable for delivering high intensity data like facial expressions. Table below presents some of the popular haptic technologies developed for communicating high bandwidth interpersonal interaction data. While these provide opportunity to appreciate the possibility of haptic information delivery, the end goal was not clearly defined towards delineating fine facial movements that are deemed important in typical dyadic social interactions. First we investigate various haptic interpersonal communication technologies that have been developed over the past several years, following which we discuss the choice of technology and any specific justifications for its use in delivering facial expressions.

6.1 Related Work in Haptic Interpersonal Communication Interfaces

Haptics has recently seen a uphill trend in terms of the number of interfaces that are being suggested by academia and the industry. It is not surprising to notice a plethora of devices that are intended for personal interpersonal communications. Tables 6.1 through 6.4 discuss various interpersonal communication technologies that were mostly developed for enriching remote interactions between humans. This list provides an appreciation for the types of technologies that have been implemented in the human-human communication enrichment space. Following these tables, we discuss the technology that we propose to use for delivering facial expression information.

Table 6.1: Haptic interpersonal interaction enrichment devices

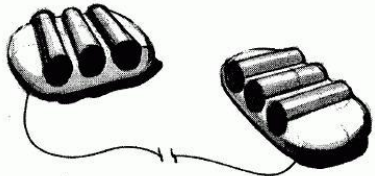
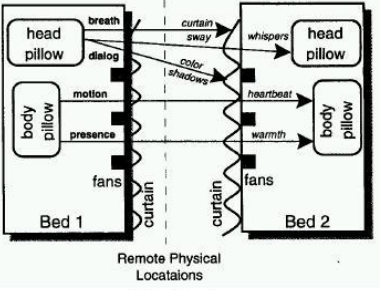
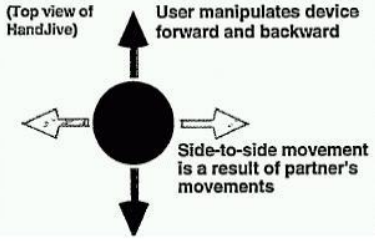
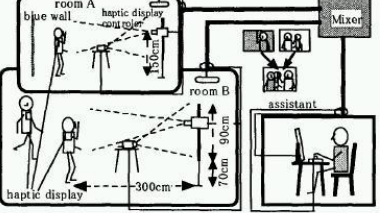
Device	Actuation	Major Organ	Non-verbal Cue	Application	User Experience	Device	Device
inTouch [253]	Vibration Pressure Texture	Hand	Social Touch	- Interpersonal communication through remote touch	No Testing		
The Bed [254]	Pressure Texture Temperature	Body	Social Touch	Remote interpersonal intimate communication Self report	System found to produce feelings of intimacy.		
HandJive [255]	Pressure Proprioception	Hand	Handshake Touch	Interpersonal communication through new cueing language	No Data		
HyperMirror [256]	Pressure	Shoulders	interpersonal distance, relative position, crossing paths	-Remote crossing of paths. -Initiating interaction in strangers across distance. Tap to initiate conversation.	Eye contact was made across distant users. Tap signal aroused attention. Crossing paths initiated conversations.		

Table 6.2: Haptic interpersonal interaction enrichment devices contd.


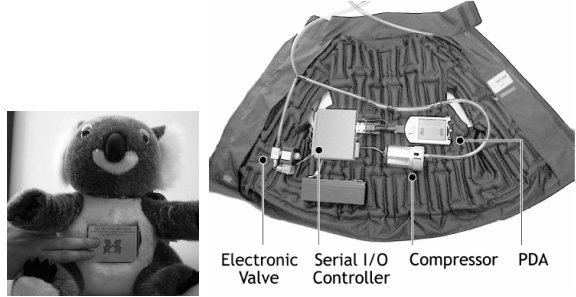
Device	Actuation	Major Organ	Non-verbal Cue	Application	User Experience	Device																												
Com Touch [257]	Vibration Pressure	Hands	Touch Emotions	Bidirectional operation. 24 subject tested. Remote participant squeezes one end and a recipient at the other end feels vibrations.	Subjects came up with their own cueing. 83% of participants used at least one gesture. 67% developed their own gestures.																													
Haptic Instant Messenger [258]	Audio Vibrations	Hands	Emotions	Based on user selections at a remote location, haptic and audio codes are transmitted to the receiver.	No user testing.	<table border="1"> <thead> <tr> <th>Icon</th> <th>Emoticon</th> <th>Meaning</th> <th>Hapticon</th> </tr> </thead> <tbody> <tr> <td></td> <td>:)</td> <td>regular smile</td> <td></td> </tr> <tr> <td></td> <td>:D</td> <td>big smile</td> <td></td> </tr> <tr> <td></td> <td>:(</td> <td>sad face</td> <td></td> </tr> <tr> <td></td> <td>;-)</td> <td>wink</td> <td></td> </tr> <tr> <td></td> <td>(k)</td> <td>kiss</td> <td></td> </tr> <tr> <td></td> <td>:\$</td> <td>embarrassed</td> <td></td> </tr> </tbody> </table>	Icon	Emoticon	Meaning	Hapticon		:)	regular smile			:D	big smile			:(sad face			;-)	wink			(k)	kiss			:\$	embarrassed	
Icon	Emoticon	Meaning	Hapticon																															
	:)	regular smile																																
	:D	big smile																																
	:(sad face																																
	;-)	wink																																
	(k)	kiss																																
	:\$	embarrassed																																
Hug over Distance [259]	Pressure Proprioception	Upper Body	Touch Hug	At one end the user rubs tummy of a stuffed toy and based on the pressure applied, air bags are filled at the remote end to simulate hug.	Air compressor at the receivers end makes a lot of noise. Six couple focus group found the concept weird.																													

Table 6.3: Haptic interpersonal interaction enrichment devices contd.



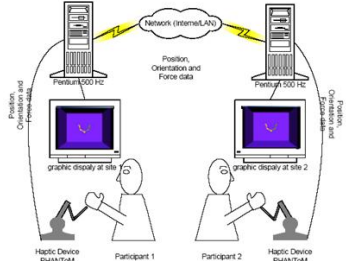


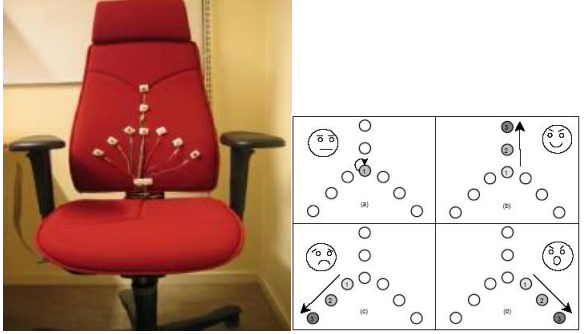
Device	Actuation	Major Organ	Non-verbal Cue	Application	User Experience	Device
VinroBod [260]	Pressure, Temperature	Hands	Touch	Convey remote interpersonal cues.	15 subjects found the device useful and intuitive	
What's Shaking [260]	Vibration Temperature	Hands	Proxemics	Heat corresponds to the number of people. Vibration corresponds to the amount of activity in the environment	12 subjects found the glove intuitive and were able to identify activity around them.	
Tele Handshake [261]	Proprioception	Hnads	Touch	Remote handshake between interaction partners	65% Satisfaction 55% Convincing 60% Intuitive	

Table 6.4: Haptic interpersonal interaction enrichment devices contd.

Device	Actuation	Major Organ	Non-verbal Cue	Application	User Experience	Device
TapTap [262]	Pressure	Shoulders	Touch Tap	Solenoids and vibrators used on the shoulder to simulate tapping.	8 men and 8 women tested on the device found based on the tap, it reminded them of someone.	
United Pulse [263]	Vibrations	Finger	Intimacy	Vibrators on the ring stimulated to initiate communication between remote couple. Simulated heart beats were delivered	20 couples tested with the device. 22 liked the idea. 5 were irritated.	 <p>heart rate monitor inside microcontroller and wireless connection to the mobile (Bluetooth, Radio Frequency)</p>
Haptic Chair [251]	Vibrations	Back	Emotions	Vibrations corresponding to emotions are delivered to the back of the user. Has sensing of the emotions inbuilt through vision technologies	3 expressions tested. 100% recognition on expressions. 10% of participants complained of cognitive load.	

6.2 Proposed Visuo-Haptic Sensory Substitution Device

In order to deliver the high bandwidth visual facial expression data, we resort to choosing a modality for communication and the human body organ where the data will be delivered to. To this end, we chose vibrotactile stimulators for encoding the facial movement patterns to temporal vibrations. Further, the facial expressions themselves were delivered as spatial pattern on a matrix of haptic actuators placed in contact with the dorsal surface of the fingers, which allows both spatial and temporal mapping of vibration patterns. The fingers have the largest tactile representation in the brain after the tongue. Together, the fingers have the largest projection on the cortical surface, similar in sensitivity to the tongue and lips (See Figure 6.1 [2]). The concentrated neuronal mapping of the fingers allow for a very high sensitivity (both spatial and temporal resolution) making them an ideal candidate for sensory substitution. Further, to allow functional operation of the user's hands, the vibratos are placed on the dorsal surface of the fingers. The following table shows a list of applications that have called for the use of haptic glove based solutions.

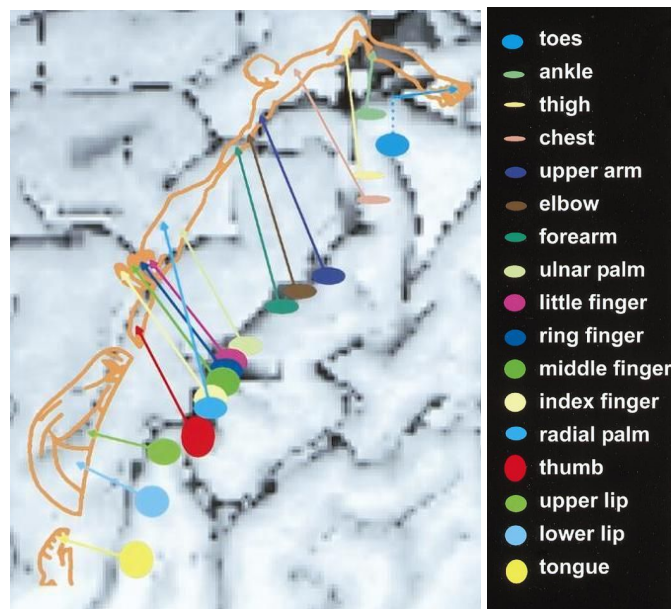


Figure 6.1: Somatosensory homonculus as mapped through magnetoencephalography [2].

Table 6.5: Related work in the development of haptic gloves for information delivery

Ref.	Application	No. of vibrators	Location of vibrators	Vibration Pattern	Encoding	Experiments	User Study
[264]	Convey color information to people who are blind.	3	Distal phalanges of index, middle and ring fingers (T). - Three phalanges of the index finger. (O)	- Continuous on all three vibrators (S). - 0.5s time gap between vibrators (D).	- Encode R, G and B channel to each of the 3 vibrators. - Amplitude of vibration proportional to the intensity of the color channel.	- Convey only colors individually (C). - Allow users to explore a down sampled color image using a mouse (I).	- 5 participants who are blind. - 2 sighted participants. - COS: 71% - CTS: 87% - ITS: 100% - IOD: 67% - IOS: 92% - COD: 87% - CTD: 90%
[265]	Vibrotactile cueing to improve target acquisition in virtual 2D environment using mouse as input.	4	-2 on the lower part of the palm just above the wrist. - 2 on the back of the lower palm just above the wrist.	- 100ms vibratory cues to indicate direction of the target and on-target signals. - Frequency of the vibration was proportional to the direction and distance from the target location.	- Two vibrators were turned on to indicate arrival on a target.	- Expt 1 tested vibrators on the front and back of palm. - Expt 2 tested continuous distance cueing with suppressing or increasing frequency as the target is approached.	- The location of the tractors did not have an effect. Front and Back worked the same. - Suppressing the frequency as the user approaches the target worked better than enhancing.
[266]	Vibrotactile array for delivering distance to an obstacle from a wheelchair driven by a visually impaired person.	9	Array on the front of the palm in a 3x3 matrix.	- Warning signals - Spatial obstacle location signal. - Direction conveyance to the user.	- Warning signal vibrates all vibrators. - Spatial location of an obstacle is sent in the particular motor with near, medium and far range to obstacle. - Direction cue vibrates the center motor with two pulses and then vibrates motor of the desired direction.	- No user testing done yet	- No user testing done yet

Table 6.6: Related work in the development of haptic gloves for information delivery contd.

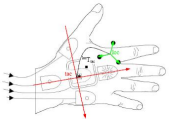
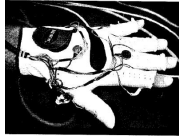
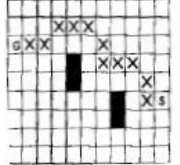
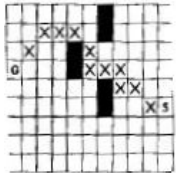
Ref.	Application	No. of vibrators	Location of vibrators	Vibration Pattern	Encoding	Experiments	User Study
[267] [268]	Vibrotactile cues for navigating surgeons hand during surgery.	4	See Figure. 	- Continuous vibrations based on the amount of off target displacement	- Optical tracking of visual markers on the surgeons hand is translated to vibrotactile cues to give off-center information.	- Subjects were required to move a surgical tool to the target location. - Subjects react to varying impulse input as required.	No quantification provided in the paper.
[268]	Field of view in front of individual who is blind is captured with a camera and translated to vibrotactile cues corresponding to a depth map.	No data	No specific information provided. 	- Magnitude of vibration is directly proportional to distance to obstacle. - Frequency of vibration is inversely proportional to the confidence in depth measurement.	The image from the camera is used to determine a depth map of obstacles in front of the user and is translated into vibrotactile cues.	Two obstacle courses were set within the laboratory environment and the participants were required to navigate the course. Course 1:  Course 2: 	- 9 participants, 3 blind and 6 with low vision. - Course 1: Traveled the minimal hitting path 65% with their existing navigation aid and increased to 75% with the glove. - Course 2: Traveled the minimal hitting path 65% with their existing navigation aid and decreased to 57% with the glove.

Table 6.7: Related work in the development of haptic gloves for information delivery contd.

Ref.	Application	No. of vibrators	Location of vibrators	Vibration Pattern	Encoding	Experiments	User Study
[269] [270]	Framework for delivering haptic data along with audio video data from an entertainment perspective. Specifically, adding a haptic layer to the MPEG 4 audio video coding.	76	Vibrotactors are added all over the glove both on top and bottom of the hand. No specific configuration pattern is discussed in the paper.	Custom designed vibration patterns that take into account all the vibrators on the glove.	Manually encoded by entertainment specialists based on the movie and the scene.	No user study.	No user study.
[271]	Using vibrators to convey slip information in a prehensile glove.	5	Fingertips of the five fingers.	Motion sensors (optical motion sensor similar to the one used in an optical mouse) mounted outside the glove on the finger tips measure the slip of an object. The slip information measured as optic flow is conveyed to the vibrator as varying frequency.	Slip motion is proportional to the frequency of vibration.	12 subjects. Users placed the glove on a surface that was laterally pulled from under the glove and the reaction time was measured by asking the participants to press a button with their free hand. Experiment was conducted with bare hands, with a prehensile glove without vibrators and with the slip glove.	<u>Mean reaction time:</u> Bare hand: 0.214s Normal Glove: 1.669s Slip Glove: 0.483s <u>Percent Failure:</u> Bare hand: 0% Normal Glove: 27.8% Slip Glove: 5.6%

These technologies have proved the viability for glove based vibrotactile sensory augmentations. Adding to these findings, the experiments described in this paper addresses the vibrotactile sensory abilities of the dorsal surface of the fingers, especially for sensory substitution.

6.3 The Vibrotactile Glove

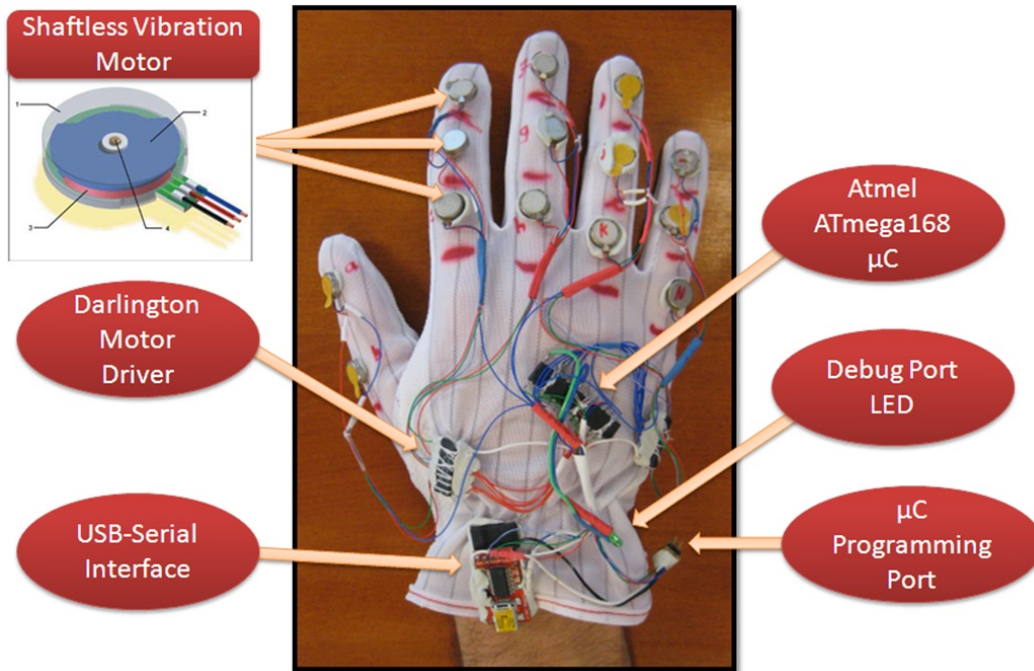


Figure 6.2: The Vibrotactile Glove.

As shown in Figure 6.2, 14 shaftless vibration motors are installed on the dorsal surface of the phalanges of a stretchable anti-static glove - corresponding to the 14 phalanges of the human hand. Each vibrator has an effective displacement of 1.5mm @ 55Hz with an effective acceleration along X, Y and Z of $X_g = 0.38g$, $Y_g = 0.29g$ and $Z_g = 1.08g$, respectively, with Z axis perpendicular to the skin. The g here refers to the acceleration due to gravity, which is equal to $9.8m/s^2$ roller via a serial port that is translated to USB for interfacing with any generic computing element. The USB port also provides the necessary power for the operation of the glove. Through these commands, the microcontroller allows precise simultaneous control of three dimensions of vibrations, namely, the intensity of vi-

bration, the location of vibration and the duration of vibration. The location of vibration is controlled by choosing the appropriate output port of the microcontroller; the duration of vibration is controlled by the onboard timer; the intensity of the vibration is controlled via simulated Pulse Width Modulation (PWM) on the output ports. In order to isolate and protect the controller from the vibration motor induced back-EMF (electromotive force) two 7-array Darlington transistors with opto-isolation are used between the controller and the motors.

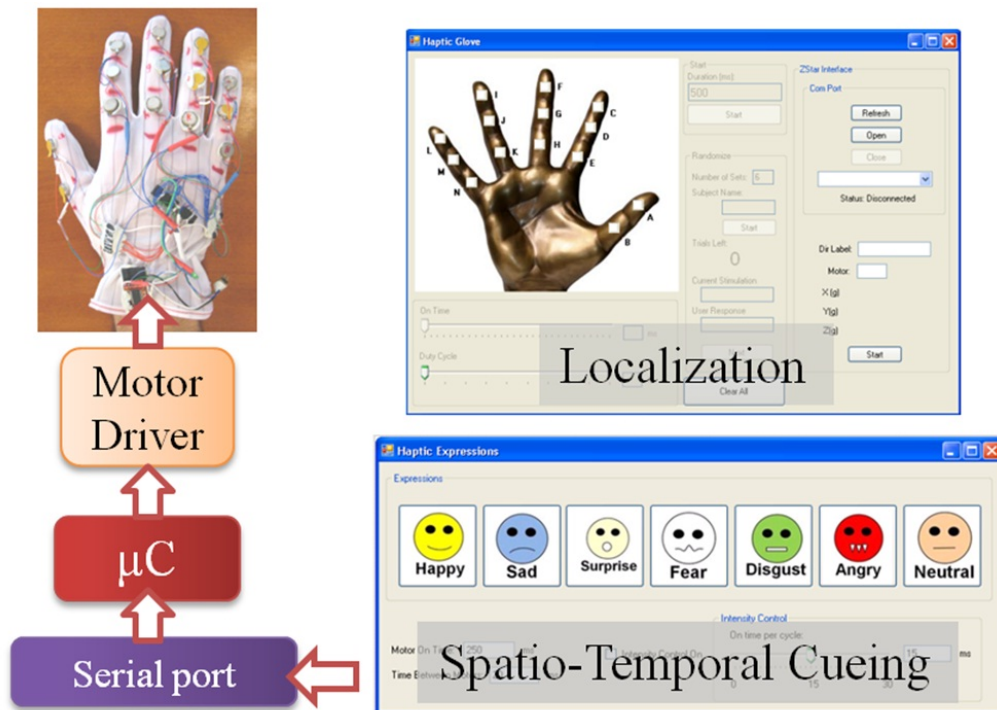


Figure 6.3: Localization and spatio-temporal cueing software used for the vibrotactile glove.

The software to control the vibrations on the glove is shown in Figure 6.3. Two independent programs were developed to explore the localization capabilities of users, and for testing the ability of users to identify spatio-temporal cues, the spatio-temporal cues here refer to the facial expressions that were being delivered from the Temporal Exemplar-based Bayesian Network decision system. The number of spatio-temporal patterns were decided based on the number of classes N that were chosen in the TEBN, in our case 6 corresponding to the 6 basic human expressions. We also added one more class, Neutral

class, to represent that the interaction partner was not displaying any expression. The serial port interface for the glove is designed similar to the popular Hayes AT command set. Activation commands are passed as ASCII strings that are interpreted by the microcontroller to activate the appropriate motor for the requested duration and intensity of vibration. The details of the cueing patterns are presented in the next section.

6.4 Haptic Cueing to Test Localization and Spatio-Temporal Mapping

6.4.1 Localization

To determine how well users were able to perceive the vibratory patterns on the phalanges, vibrators were excited at randomly selected locations. Each excitation was applied at 100% intensity and duration of 5 seconds. The localization experiments were focused on studying the vibrotactile detection capabilities of the individual phalanges, fingers as a whole, and groupings based on the distance of the phalanges from the palm (distant, intermediate and proximal phalange) as shown in Figure 6.4.

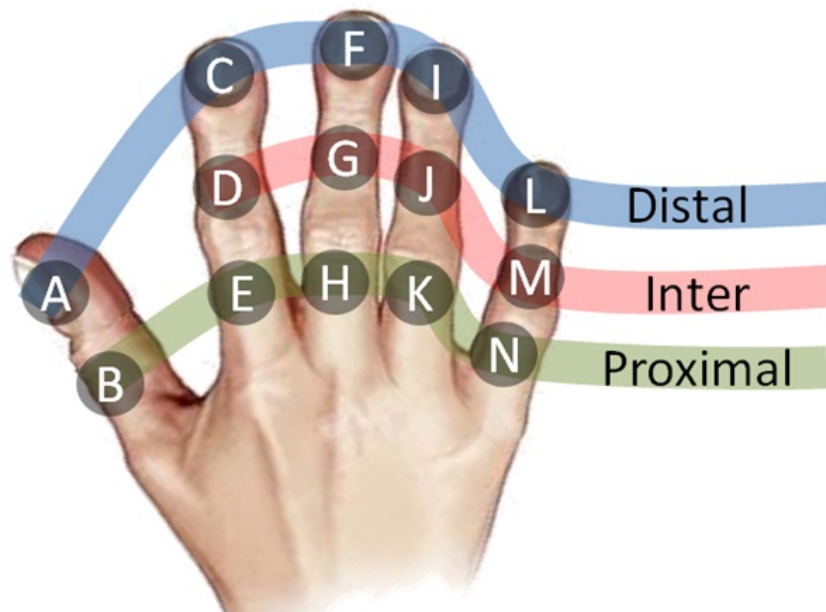


Figure 6.4: Phalange naming convention and grouping based on the anatomical distances.

6.4.2 Spatio-Temporal Cueing

While the versatility of the VibroGlove allows it to be used for various applications, here we discuss the specific application of delivering the six basic facial expressions, along with the neutral face, of an interaction partner to a user who is visually disabled. Humans rely heavily on the shape of the mouth and the eye area to decipher facial expressions. Motivated from this, we focused only on the mouth area to design spatio-temporal haptic alternates for facial expressions. We used only the three central fingers on the glove: 9 vibrators, as shown in Figure 4. In order to represent the seven facial expressions, we designed haptic expression icons that were motivated by two important factors: 1) Icons similar to the visual emoticon that are already in popular use, like Happy, Sad, Surprise and Neutral, where the mouth shapes prominently represent the expression, and 2) Icons like Anger, Fear and Disgust where the mouth area alone does not convey the expression, thereby forcing us to create haptic icons that could evoke a sense of the expression in question. Figure 6.5 provides details of the haptic expression icons. All 7 patterns were designed to be 750ms long with each motor vibrating for at least 50ms. These numbers were determined based on pilot studies where we found that participants could not isolate vibrations if the duration was less than 50ms long. Further, patterns longer than 800ms were considered to be too long by the participants, while patterns shorter than 600 ms were confusing, and training phase accuracies were unacceptable.

6.4.2.1 Group 1 - The visual emoticon motivated haptic icons:

The Group 1 haptic expression icons primarily represent popular emoticons that are in wide use within the Instant Messaging community. These icons mostly model the shape of the mouth.

- 1) Happy is represented by a U shaped pattern,
- 2) Sad by an inverted U,

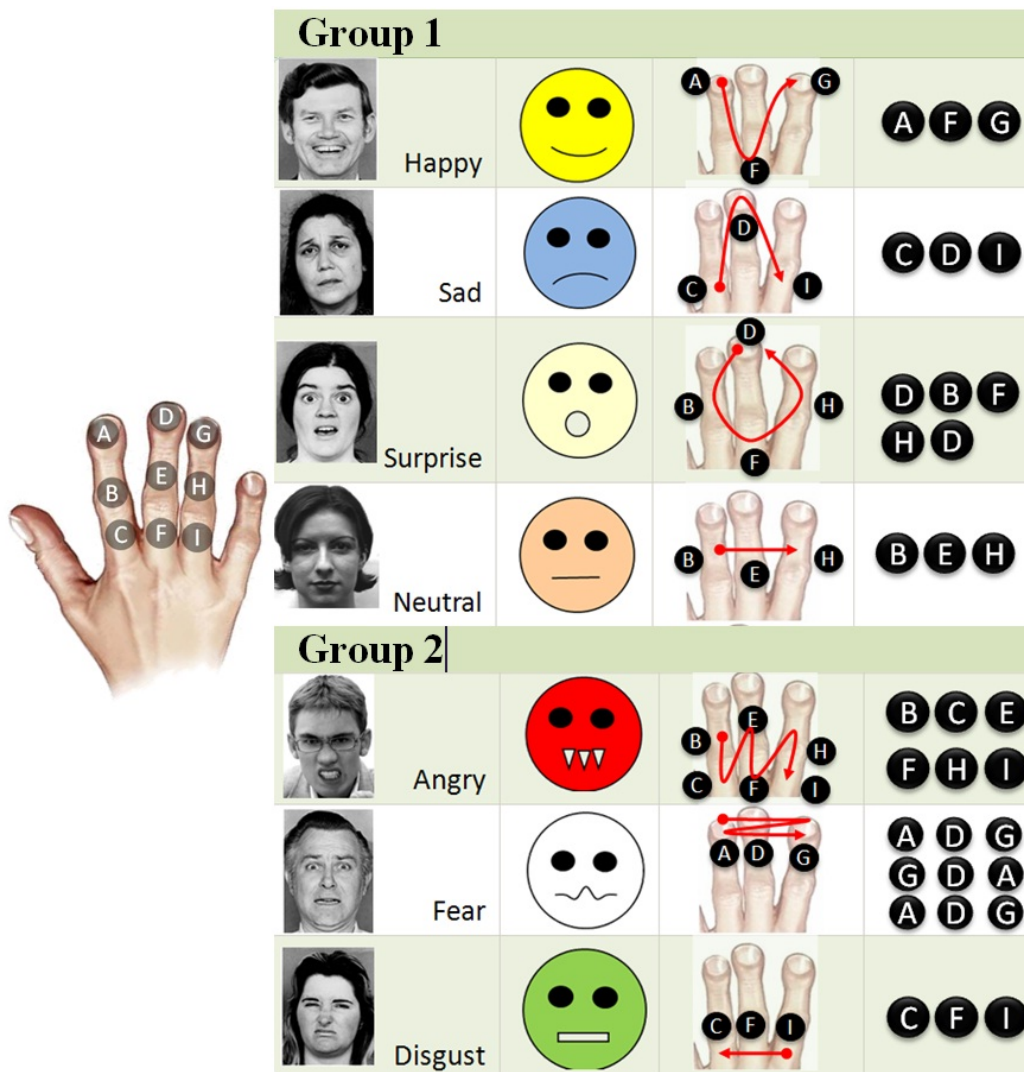


Figure 6.5: Mapping of Group 1 and Group 2 haptic expression icons to the central three fingers (9 Phalanges) of the vibrotactile glove. In the expression mapping chart, Columns 1 to 3 represent the expression. Column 4 shows the spatial mapping of vibrations. Column 5 shows the temporal mapping of the vibrations.

- 3) Surprise by a circle ○, and
- 4) Neutral by a straight line —.

6.4.2.2 Group 2 - The auxiliary haptic icons:

Anger, Fear and Disgust cannot be conveyed through the appearance of mouth alone. To this end, we resorted to defining haptic patterns that were unique from what was already defined for Group 1, while keeping in mind a need to represent the underlying expression in question.

- 1) Anger is represented by successive vibrations on six lower phalanges representing an open mouth showing teeth during an expression of anger;
- 2) Fear is represented by very brief vibrations on the dorsal phalanges of the central 3 fingers in three quick successive vibration sequences representing a fast emotional response that people show towards fear, and
- 3) Disgust is represented through a vibration pattern going from right to left on the bottom phalanges of the central fingers corresponding to a slightly opened mouth during the display of disgust.

6.5 Research Hypothesis

6.5.1 *Localization*

While testing the localization capabilities of the haptic glove, three distinct and correlated hypotheses were tested. These hypotheses are related to the individual phalange localization, localization per finger, and localization on the phalange groups based on their distance from the palm.

6.5.1.1 Hypothesis 1:

- a) The recognition rates per phalange will be above chance (50%);

- b) The mean recognition rate per phalange will not be significantly different between any two phalanges.

6.5.1.2 Hypothesis 2:

- a) The recognition rates per finger will be above chance (50%);
- b) The mean recognition rate per finger will not be significantly different between two fingers.

6.5.1.3 Hypothesis 3:

- a) The recognition rates per phalange group (distal, intermediate, proximal) will be above chance (50%);
- b) The mean recognition rate per phalange group will not be significantly different between two phalange groups.

6.5.2 *Spatio-Temporal Cueing*

Similar to the localization experiments, the hypotheses relating to the spatio-temporal cueing relates to the ability of the users to recognize the individual expression and also the two groups of expressions as identified in Section IV B.

6.5.2.1 Hypothesis 4:

- a) The recognition rates for the spatio-temporal expression patterns will be above chance (50%);
- b) The mean recognition rate per expression will not be significantly different between any two expressions.

6.5.2.2 Hypothesis 5:

- a) The recognition rates per expression group (Group 1 and 2) will be above chance (50%);
- b) The mean recognition rates between the two groups will not be significantly different.

6.6 Experiments and Analysis Methodology

Two independent and consecutive experiments were conducted to test the localization and spatio-temporal cue identification capabilities of users. Participants were engaged for the entire time of the two experiments and the localization experiments preceded the spatio-temporal cueing experiments. In the spatio-temporal experiments, along with the accuracy of recognizing the spatio-temporal vibrotactile cues, we were also interested in knowing how quickly the participants were able to recognize the expressions. The duration for recognition is very important in social interactions as the human face changes drastically over short time. Experiments have shown that expressions vary anywhere from 1 to 5 seconds ([5], Page 322). Conforming to these time scales, it is important that any device developed towards enriching social experience should react in real-social-time towards facilitating smooth interpersonal interaction.

6.6.1 *Participants*

The experiments were conducted with one individual who was blind and 11 other participants who were sighted but blindfolded during the experiment. It is important to note that the individual who was blind had lost his sight after 25 years of having vision. To a large extent, this individual could correlate with the Group 1 haptic expressions, of the spatio-temporal cueing experiment, to his visual experiences from the past. None of the participants had any obvious medical conditions that prevented them from perceiving the vibrotactile stimulations on their right hand.

6.6.2 *Procedure*

Once the subjects wore the glove, they were seated in a chair with a blindfold and asked to keep their hand on their lap in the most comfortable position. Both the localization and spatio-temporal cueing experiments were conducted in three successive phases, namely, Familiarization phase, Training phase and Testing phase.

Familiarization Phase: Subjects were first familiarized with the various vibration patterns by presenting them in order - each phalange for the localization experiment and each facial expression for the spatio-temporal experiments. During this phase, the corresponding location or the facial expression was spoken aloud by the experimenter. The familiarization was continued until the subjects were comfortable in remembering all the locations and expressions.

Training Phase: The familiarization was followed by the training phase in which all the fourteen vibration locations and seven facial expression patterns were presented in random order, in multiple sets, and subjects were asked to identify them by speaking them out. The experimenter confirmed any correct response, and corrected incorrect responses. Subjects had to demonstrate 100% recognition on at least one set of all fourteen locations and seven expressions before moving to the testing phase. Note that, the speaking was replaced by the use of a keyboard in the spatio-temporal experiments where the participants were asked to type their answers directly into a keypad having 7 keys corresponding to the 6 basic expressions and the neutral face. The keypad allowed us to capture the exact time taken by the participants to arrive at the decision and respond with a key press. A 15 minute time limit was placed on the training irrespective of the training accuracy.

Testing Phase: The testing phase was similar to the training phase except the experimenter did not provide feedback to subjects, and each location and expression pattern was randomly presented 10 times making a total of 14 locations x 10 trials = 140 localization results, 7 expressions x 10 trials = 70 expression results. The subjects were given 5 seconds per trial to respond.

6.6.3 Analysis

In order to test the hypotheses presented in Section 6.5 (relating to the localization and spatio-temporal cueing experiments), three related analyses were carried out, namely, a) location or expression recognition rate, b) One-way analysis of variance (ANOVA) on the

recognition accuracies, and c) Tuckey Honestly Significant Difference (HSD) test to determine the mutual performance of location and expression results. The details of the three techniques are discussed below.

6.6.3.1 Recognition Accuracies

As the name suggests, the recognition accuracies measure the average true positive rate of recognition on the localization and expression recognition experiments. Along with the mean recognition rate, the deviation in the recognition rates across the 11 participants is also shown.

6.6.3.2 ANOVA

One-way analysis of variance is a statistical tool used for comparing two or more sample groups to test null hypothesis that the samples were drawn from different populations using an F-distribution. The F statistic is derived from the sample means and the group means. Following the central limit theorem, if the samples are drawn from the same population, the variance between group means have to be smaller than the sample variance. A higher ratio of the variances justifies the null hypothesis, else it's rejected. The results of ANOVA are reported as p-value scores from the F-statistic with the dimensions $(k-1)$ and $(n-1)$, where the k is the number of groups and n is the number of samples. Lower the p-value higher is the chance of accepting the null hypothesis and vice versa.

6.6.3.3 Tuckey HSD

While ANOVA tests for a chance that the samples could have been derived from different populations, the Tuckey Honestly Significant Difference (HSD) test relies only on the group means to determine if there is a significant difference between groups of samples. A significant difference calls for reasoning to suspect/explain performance differences within groups derived from a single sample set. Unlike ANONA, where all the groups and samples are combined into the dimensions of comparison, HSD allows individual group-wise comparisons, providing for an opportunity to identify which groups are performing differently

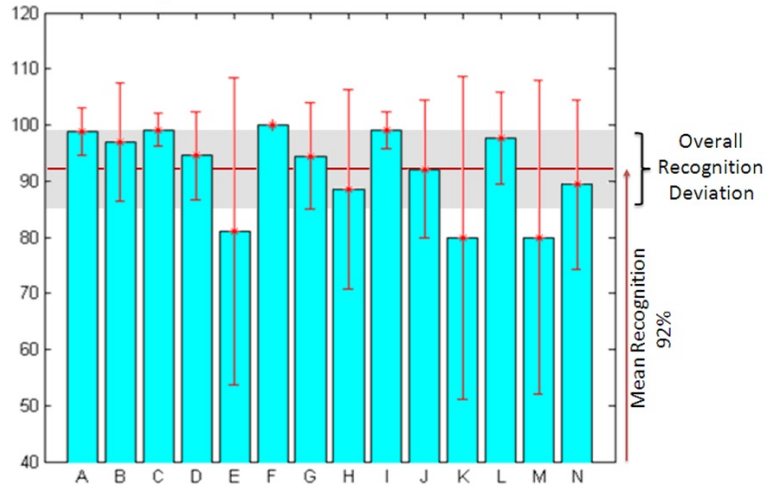
from others. Mostly reported as a ratio of the group mean difference and the standard group mean error, it is possible to quickly identify the significant group differences. In the results section below, the group means and the standard errors are plotted as circles and whiskers, respectively. Significant difference is established when one group's standard error stretch is beyond the scope of any of the other groups.

6.7 Results of the Experiments

In this section, five sets of results are presented. Each set presents the three analyses that were described in Section 6.6.3. Three of the five sets correspond to the localization experiments, while the other two correspond to the spatio-temporal experiments.

6.7.1 Localization Experiments

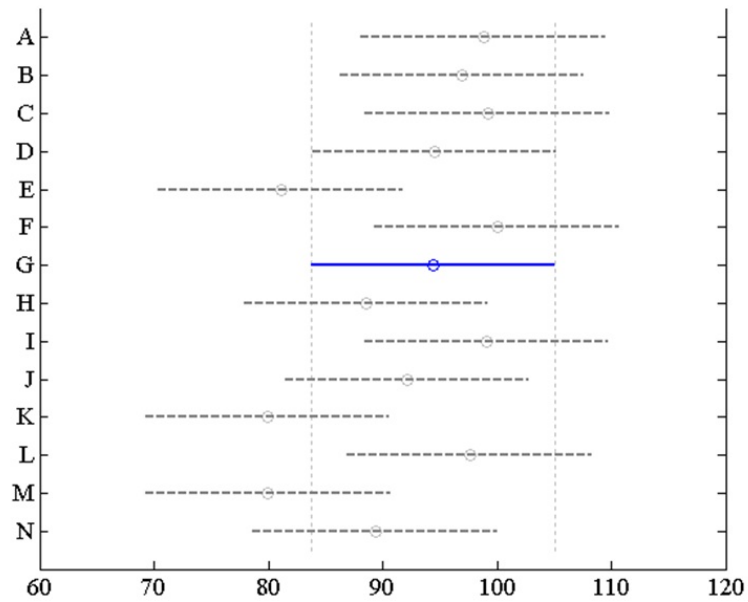
6.7.1.1 Phalange Level Localization



(a)

	SSE	DoF	Mean Sq.	F	p-value
Between Groups	8488.7	13	652.97	2.68	0.002
Within Groups	37529.6	154	243.69		
Total	46018.2	167			

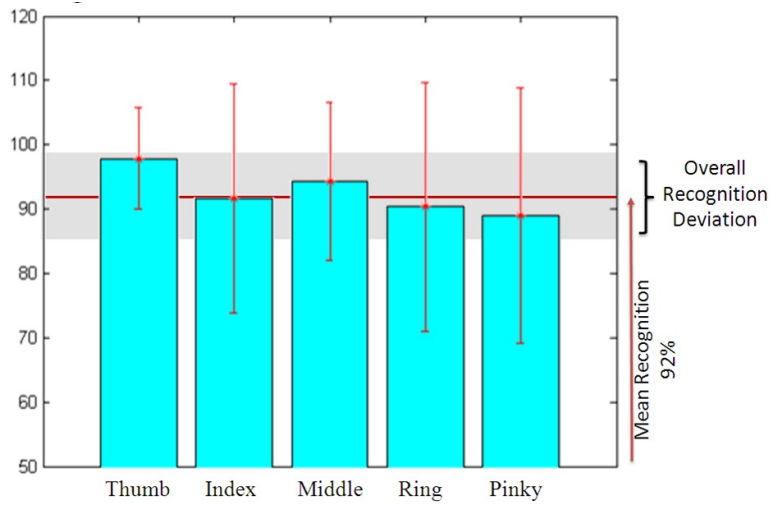
(b)



(c)

Figure 6.6: (a) Recognition Accuracies; (b) ANOVA; (c) HSD

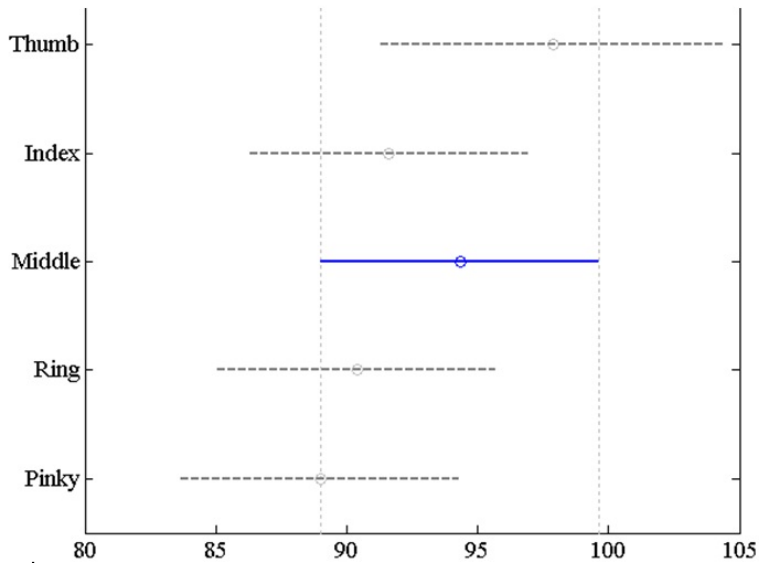
6.7.1.2 Finger Level Localization



(a)

	SSE	DoF	Mean Sq.	F	p-value
Between Groups	1440.9	4	360.234	1.32	0.2657
Within Groups	44577.3	163	273.48		
Total	46018.2	167			

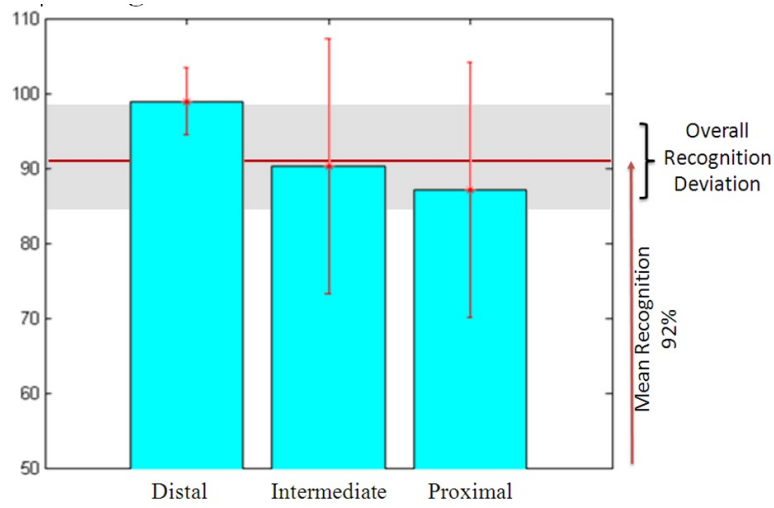
(b)



(c)

Figure 6.7: (a) Recognition Accuracies; (b) ANOVA; (c) HSD

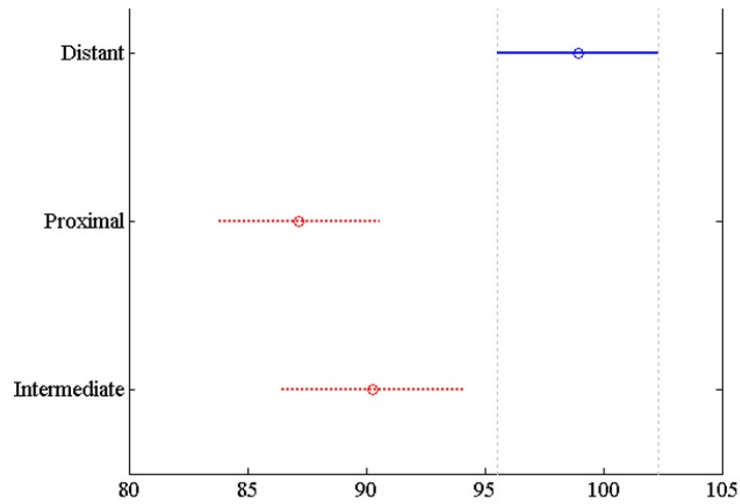
6.7.1.3 Phalange Position Localization



(a)

	SSE	DoF	Mean Sq.	F	p-value
Between Groups	4408.7	2	2204.33	8.74	0.0002
Within Groups	41609.6	165	252.18		
Total	46018.2	167			

(b)

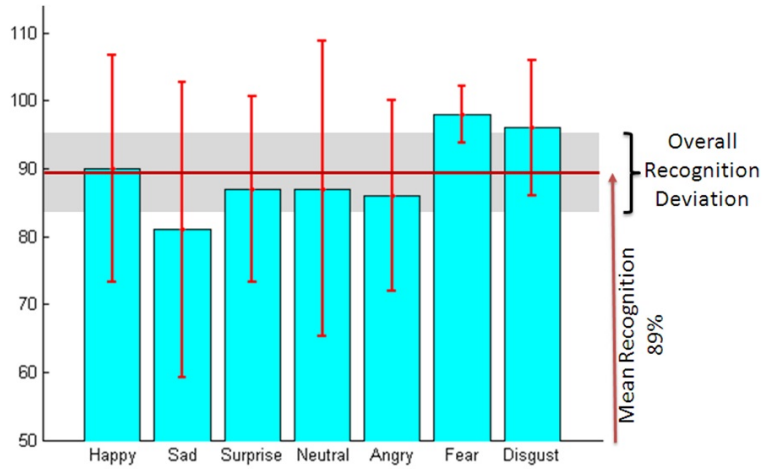


(c)

Figure 6.8: (a) Recognition Accuracies; (b) ANOVA; (c) HSD

6.7.2 Spatio-Temporal Experiments

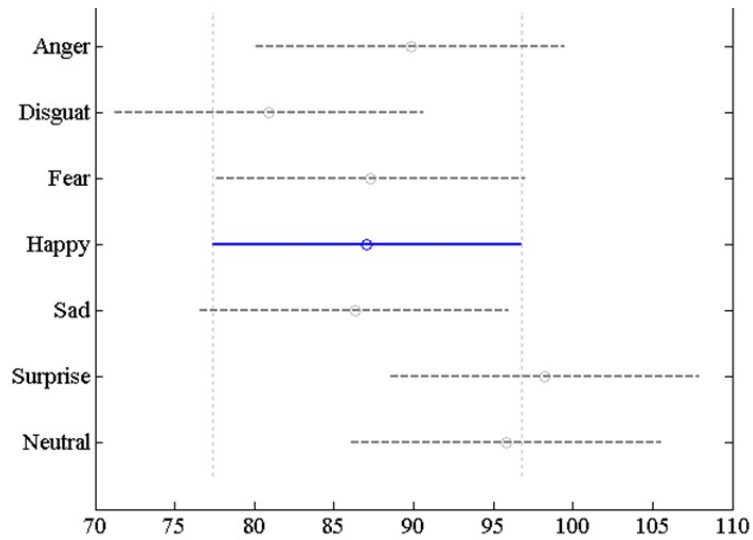
6.7.2.1 Individual Spatio-temporal Patterns



(a)

	SSE	DoF	Mean Sq.	F	p-value
Between Groups	2530.4	6	421.72	1.71	0.1299
Within Groups	18990.9	77	246.63		
Total	21521.3	83			

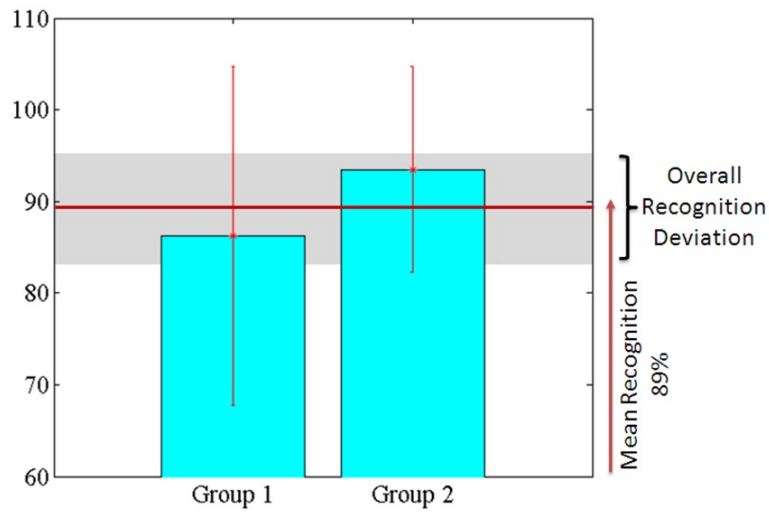
(b)



(c)

Figure 6.9: (a) Recognition Accuracies; (b) ANOVA; (c) HSD

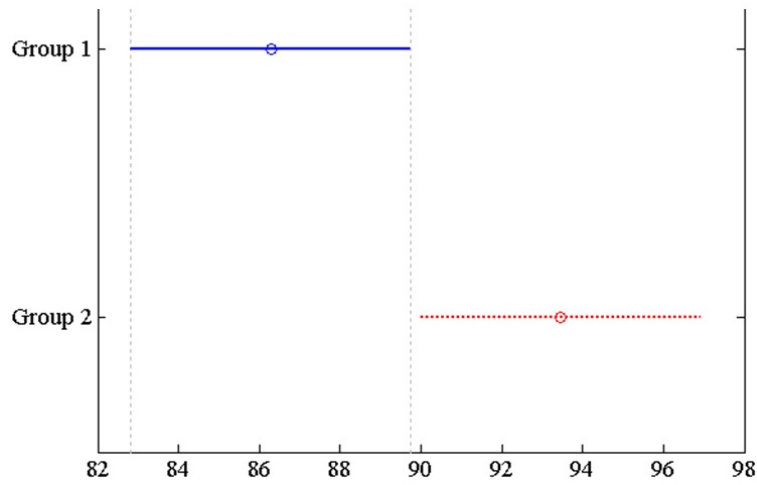
6.7.2.2 Comparison of Spatio-temporal Cueing groups



(a)

	SSE	DoF	Mean Sq.	F	p-value
Between Groups	1058.2	1	1058.16	4.24	0.0426
Within Groups	20463.1	82	249.55		
Total	21521.3	83			

(b)



(c)

Figure 6.10: (a) Recognition Accuracies; (b) ANOVA; (c) HSD

From Figure 6.6 through 6.11, it can be seen that the part (a) of the five research hypotheses can be answered immediately. Users of the device found it convenient to localize vibration patterns and identify expressions easily. The localization recognition accuracy was at 92% (SD: 7.5%), while the expression recognition rate was measured at 89% (SD: 5.9%). This validates the null hypothesis that the users were able to localize and identify vibrotactile patterns well above average.

Investigating part (b) of the five hypotheses reveals interesting insights into the user's abilities to detect and localize vibrotactile stimulations. From Figure 6.6(b) and 6.6(c) it can be concluded that Hypothesis 1(b) is accepted; both ANOVA HSD tests reveal no significant difference between phalange performances. Similarly from Figure 6.7(b) and 6.7(c), it can be concluded that there is no significant difference in the performance between fingers. Figure 6.9(b) and 6.9(c) accepts the Hypothesis 4(b) and we see no significant difference in the mean performance of the seven spatio-temporal cueing patterns of facial expressions.

In contrast, from Figure 6.8(b) and 6.8(c), we see that the Hypothesis 3(b) is rejected as user performance diminished at the proximal phalanges. This could be attributed to the fact that the vibration motors are very closely placed next to one another at the proximal phalanges which may cause inter-motor vibrations. From Figure 6.11(b) and 6.11(c), we see that Group 2 performance was much higher than Group 1 rejecting the Hypothesis 5(b). Studies are underway to determine the nature of the haptic cues in Group 2 that make them significantly better than Group 1. This could have been due to the fact that Group 2 cues were designed based on extensive user feedback when compared to Group 1 expressions which were designed based on popular visual emoticons.

Figure 6.7 shows the confusion matrix for all the seven expressions. The diagonals correspond to the bar graph shown in Figure 6.6. The off-diagonal elements represent the confusion between expressions. These off-diagonal elements provide insight into the parameters that control effective and responsive haptic patterns. While subjects confused Sad and Neutral expressions with various others (mostly in Group 1), Anger and Surprise

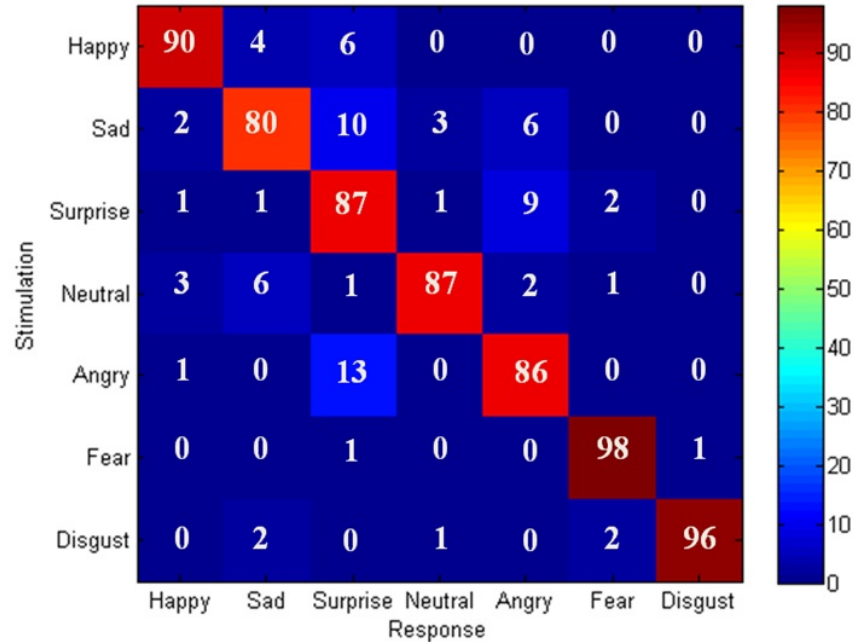


Figure 6.11: Confusion Matrix across the 12 participants. The rows are the stimulation and the columns are the responses of the participants. Each row adds to 100% (rounding error of 1%).

show exchangeability, where there is strong confusion between each other. Fear and Disgust are strongly isolated from the rest of the expressions as they were very well recognized by the subjects.

Figure 6.8 shows the average recognition performance and the average time of response for the subject who is blind. The individual was able to recognize most of the expressions at 100%, over the 70 trails.

6.7.2.3 Time for Recognition:

Figure 6.9 shows the average time taken by the subjects per expression when they recognized the haptic patterns correctly (cyan), and when they misclassified them (red). The bar graph shows excess or shortage of response time around the mean value. It can be seen that correct identification happened in just over a second (1.4s). When the subjects were not sure of the haptic pattern, they took more time to respond. This can be seen from the inverse correlation of the response time and recognition rates in Figure 6.6. The pattern for

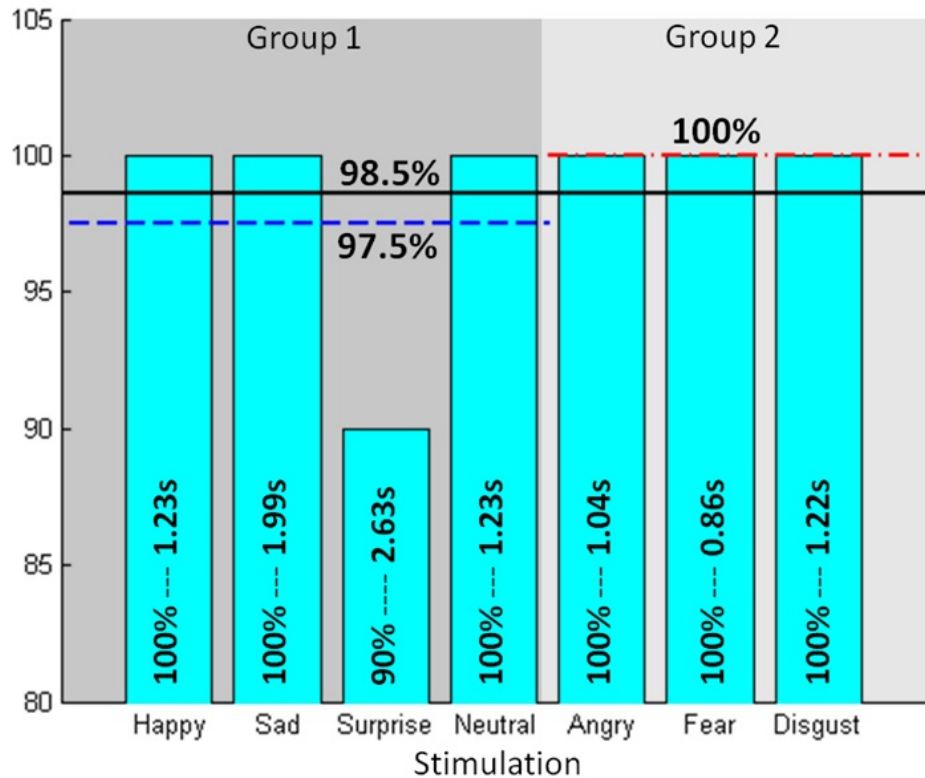


Figure 6.12: Average recognition rate and response time for the subject who is blind, for over 70 trails.

Sad had the worst performance of 81% and the corresponding response time was the highest (2s). Pattern for Fear had the best performance (98%) and least response time (765ms). This analysis can be extended to the Group level where Group 1 has a higher recognition time when compared to Group 2. Whenever the subjects responded wrong, they seem to take more time, as seen by the average incorrect response time of 2.31s (red), almost a second more than the response time for correct responses. We could not find any significant relevance between the response time for incorrect answers and the recognition rate graph. We conclude that subjects were responding with random answers once they crossed a self imagined time limit less than the 5 seconds that was provided.

6.8 Conveying Facial Expressions through Dyadic Social Situational Assistant

In this chapter we demonstrated a novel interface, a vibrotactile glove, for delivering seven facial expressions as part of dyadic social situational assistant. We have explored only

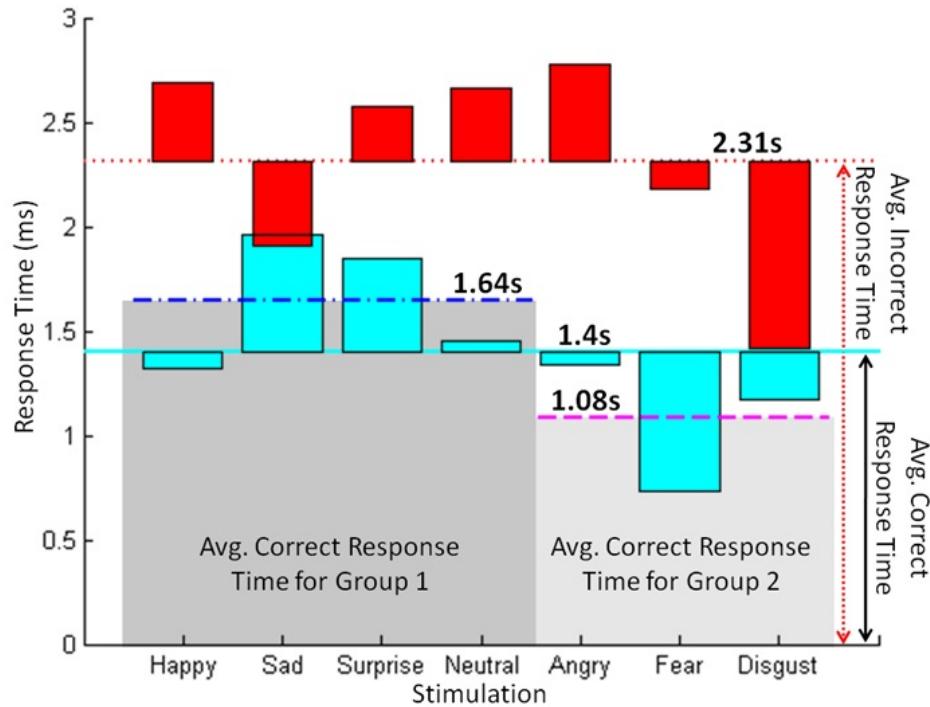


Figure 6.13: Average response time for all 12 participants. Four important results are shown above, 1) Avg. correct response time per expression (Cyan), 2) Avg. incorrect response time per expression (Red), 3) Avg. correct response time for Group 1 (Blue), and 4) Avg. correct response time for Group 2 (Magenta).

the seven classes based on the output of the Temporal Exemplar-based Bayesian Network (TEBN). If the number of classes could be increased (possibly to transmit more detailed information about the face, like Action Units of the face), the same could be done on the glove to incorporate the extension of the number of spatio-temporal classes. While the seven spatio-temporal mappings were chosen here based on their visual emoticon (Group 1) and based on the emotion that they communicated (Group 2), but if the number of classes are extended, very careful design choices would need to be made on how to extend the spatio-temporal classes without causing a sensory overload on the participants.

In the next chapter we focus on the problem of recognizing the interaction partners based on facial biometrics. We highlight the importance of face recognition in the social context.

BIOMETRICS IN SOCIAL CONTEXT: IDENTIFYING INTERACTION PARTNERS

Social interactions and social communication dynamics are very specific to the interacting partners. Social interactions of every individual with another varies depending on how they are related on the social fabric. As identified by the individuals who are blind on the online survey of important non-verbal cues (Chapter 3, identity of the interaction partner was considered as being important. In our discussions with the user population, we realized that the identity of individuals is very important while being in group discussions as they would like to know who is located where and what a certain individuals facial expressions are.

Camera-based face recognition has the potential for recognizing people at a distance, without their explicit involvement. More recently, as a result of threats to public safety, some public places (such as in Glasgow and London) have been heavily populated with video surveillance cameras. On average, a person moving through London is captured on video over 5 times a day. This offers an unprecedented basis for developing and testing face recognition as a biometric for security and surveillance.

Given this great potential, it is not surprising that many private corporations have attempted to develop and deploy face recognition systems, as an adjunct to existing video security and surveillance systems. However, the performance of these systems has been disappointing. One of the most difficult problems that face recognition researchers encounter in surveillance applications is that face databases of miscreants typically contain only frontal and profile views of each person's face, with no intermediate views. Surveillance videos captured of the same person with the same camera in the same lighting conditions might have face images that look quite different, due to pose angle variations, making it very difficult to compare captured face images to those in a database. Combine this problem with the fact that miscreants are highly motivated to disguise their identity, and the fact that face databases often contains thousands of faces, and the problem seems insurmount-

able.

Given all of these complicating factors, it is premature to rely upon face recognition systems for detecting miscreants in public places. On the other hand, the use of face recognition in controlled access applications (where users are highly motivated to cooperate, and where face database images can be both captured and tested with the same camera under the same illumination conditions) is certainly within the limitations of current face recognition algorithms.

7.0.1 Employing face recognition to facilitate social interactions

However, the application of face recognition within the proposed social situational awareness context is moderately challenging, but still potentially within the realm of possibility. This problem is simplified considerably by the fact that, on a day-to-day basis most people encounter a limited number of people whom they need to recognize. It is further simplified by the fact that people typically don't attempt to disguise their appearance in social situations. When a new person is encountered, the system could employ face detection to extract and save a sequence of face images captured during a conversation. This would provide a wide variety of facial expressions and pose angles, that could be stored in a database, and used for training a face recognition algorithm.

As people use such an assistive device over an extended period of time, they will learn both its abilities and its limitations. Conjectural information from the system can then be combined with the user's other sensory abilities (especially hearing) to jointly ascertain the identity of the person. This synergy between the user and the system relaxes some of the stringent requirements normally placed on face recognition systems.

However, such an assistive technology application still poses some significant challenges for researchers. One problem is the extreme variety of in lighting conditions encountered during normal daily activities. While there are standards for indoor office lighting that tend to provide diffuse and adequate lighting, lighting in other public places might vary considerably. For example, large windows can significantly alter lighting conditions, and

incandescent lighting is much more yellow than florescent lighting. Outdoor lighting can be quite harsh in full sunlight, and much more blue and diffuse in shadows. A person who is blind might not be aware of extreme lighting conditions, so the system would need to either (1) be tolerant of extreme variations or (2) recruit the user to ameliorate those extreme conditions.

In summary, the development of an assistive face recognition system for people who are blind provides a more tractable problem for face recognition researchers than security and surveillance applications. It imposes a somewhat less stringent set of requirements because (1) the number of people to be recognized is generally smaller, (2) facial disguise is not a serious concern, (3) multiple pose angles and facial expressions of a person can be captured as training images, and (4) the person recognition process can be a collaborative process between the system and the user.

In an attempt to provide such an assistive face recognition system, we have developed a new methodology for face recognition that detects and extracts unique features on a person's face, and then uses those features to recognize that person. Contrast this with conventional face recognition algorithms that might avoid the use of a few distinguishing features because that approach might make the system very vulnerable to disguise.

7.1 Face Recognition in Humans

For decades, scientists in various research areas have studied how humans recognize faces. Developmental psychologists have studied how human infants start to recognize faces, cognitive psychologists have studied how adolescents and adults perform face recognition; neuroscientists have studied the visual pathways and cortical regions used for recognizing faces, and neuropsychologists have attempted to integrate knowledge from neurobiological studies with face recognition research. Computer vision researchers are relatively new to this area, and have attempted to develop face recognition algorithms using image processing methods. Only recently have computer vision researchers been motivated to better understand the process by which humans recognize faces, in order to use that knowledge to

develop robust computational models. Their new interest has led to more inter-disciplinary face recognition research, which will likely aid our understanding of face recognition.

New studies have shown that humans, to a large extent, rely on both the featural and configural information in face images to recognize faces [272]. Featural information provides details about the various facial features, such as the shape and size of the nose, the eyes, and the chin. Configural information defines the locations of the facial features, with respect to each other. Psychologists Vicki Bruce and Andrew Young [273] agree with this dual representation, saying that humans create a view-centric description of a human face by relying upon feature-by-feature perceptual input, which is then combined into a structural model of the face.

Sadar et al. [274] showed that characteristic facial features are important for recognizing famous faces. For example, when they erased eye-brows from famous people's faces, face recognition by human participants was adversely affected. Young [275] showed that human participants were confused when asked to recognize faces that combined facial features from different famous faces. These studies suggest that the details of facial features are important in the recognition of faces.

However, [276] showed that the relative locations of the facial features was also very important for the recognition of faces. They collected face images of famous personalities, and then changed the aspect ratio of those images, such that the height was greatly compressed, while the width was emphasized. Surprisingly, all the resulting face images were still recognizable, despite their contorted appearance, as long as the relative locations of the features were maintained within the distorted image. This study suggests that humans can flexibly use the configural information when recognizing faces.

Another important area of research in the human perception of faces has been in understanding the medical condition of face blindness, called *prosopagnosia*. People with prosopagnosia are unable to recognize faces including their own. Until recently it was assumed that prosopagnosia was acquired often as a result of a localized stroke. However new evidence suggests that a substantial portion of the general population have a congenital

form of prosopagnosia [277]. Kennerknecht et al. [278] conducted a survey of 789 students in 2006 which showed that 17 (2.5%) suffered from congenital prosopagnosia. These students went about their daily life without realizing their disorder in face recognition.

Other studies at the Perception research centers at Harvard and Univ College of London have shown that prosopagnosics recognize people using unique personal characteristics, such as hair style, gait, clothing, and voice. These findings suggest that the detection of unique personal characteristics might provide a basis for face recognition systems to better recognize people. Since current methods of face recognition have met with only limited success, it makes sense to explore the use of this alternative approach.

Research in Own-Race Bias (ORB) in face recognition [279] has also revealed some interesting results regarding human face recognition capabilities. David Turk et al. found that, when humans are presented with new objects or new faces, they initially learn to recognize those objects and faces based on their distinctive features. Then, as familiarity increases, they incorporate configural information, moving towards holistic recognition. This study suggests that distinctive features are important during the initial stages of face recognition, and that configural information subsequently provides additional useful information.

Distinctive facial features can take many different forms. For example, after a first encounter with a person who has a handlebar moustache, we readily recognize that person by the presence of his distinctive feature. Similarly, a person with a large black mole on her face will be remembered by first-time acquaintances by that feature. Given the current limited understanding of how humans recognize faces, it makes sense to use these observations as the basis for a new approach to face recognition.

The research described in this chapter is based on the approach of identifying distinctive facial features that can be used to distinguish each person's face from other faces in a face database. In recognition of the role played by configural information in the later stages of face recognition, it also takes into account the location of these features with respect to each other. The results of our research suggest that this approach can be very

effective for distinguishing one person's face from other faces.

7.2 Our Approach to Face Recognition

Having introduced the potential for using characteristic person-specific features for face recognition, we now turn our attention towards the development of a method for discovering such features, and for using them to index face images. Then we propose a novel methodology for face recognition, using person-specific feature extraction and representation. For each person in a face database, a learning algorithm discovers a set of distinguishing features (each feature consisting of a unique local image characteristic, and a corresponding face location) that are unique to that person. This set of characteristic facial features can then be compared to the normalized face image of any person, to determine the presence or absence of those features. Because a unique set of features is used to identify each person in the database, this method effectively employs a different feature space for each person, unlike other face recognition algorithms that assign all of the face images in the database to a locality in a shared feature space. Face recognition is then accomplished by a sequence of steps, in which query face images is mapped into a locality within the feature space of each person in the database, and its position is compared to the cluster of points in that space that represents that person. The feature space in which the query face images are closest to the cluster is used to identify the query face images.

Having introduced the conceptual theory behind a person-specific characteristic feature extraction approach to face recognition, we now propose in the subsequent sections a method for detecting and extracting such features from face images, and for constructing a feature space that is unique to each person in the database.

7.3 Feature Extractors

The task of face recognition is inherently a multi-class classification problem. For every face image X , there is an associated label y that is the name of the class, i.e. the name of the person depicted in the image. While X represents the image of the person, there is no inherent constraint on whether the image is a color RGB, HUV or YCbCr image, or a

gray-scale image with a gray-scale range of 0 to 255, or even spectral representation that is extracted from the face image using Fourier transform or Wavelets. Irrespective of the image representation, the basis vectors spanning that representation are called features. The feature space spanned by these basis vectors is partitioned by the decision boundaries that ultimately define the different classes in the multi-class problem of face recognition. In this work, we choose a particular set Gabor filters as feature detectors, and each of those feature detectors for each person in the database, and that set of Gabor filters spans a unique feature space for that person.

7.3.1 Gabor Features

Gabor filters are a family of functions (sometimes called Gabor Wavelets) that are derived from a mother kernel (a Gabor Function) by varying the parameters of the kernel. As with any wavelet filters, the Gabor filters extract local spatial frequency content from the underlying image. Gabor Filters specifically capture the spatial location and spatial orientation of the intensity variations in the image underneath the filter's location. By varying the spatial frequency and the spatial scope of the filters, it is possible to extract a Gabor coefficient that partially describes the nature of the image underneath it. The coefficients obtained by filtering a locality in a face image with a set of different Gabor Filters are called Gabor Features.

7.3.1.1 Use of Gabor Filters in Face Recognition

Gabor filters have been widely used to represent the receptive field sensitivity of simple cell feature detectors in the human primary visual cortex. Recognizing this fact, Gabor features have been widely used by face recognition researchers. Over the last few years, the extensive use of Gabor wavelets as generators of feature spaces for face recognition, has led to objective studies of the strength of Gabor features for this application. For example, Shan et. al. [Shan2004] reviewed the strength of Gabor features for face recognition using an evaluation method that combined both alignment precision and recognition accuracy. Their experiments confirmed that Gabor features are robust to image variations caused by

the imprecision of facial feature localization. As indicated by Gkberk et. al. [280], several studies have concentrated on examining the importance of the Gabor kernel parameters for face analysis. These include: the weighting of Gabor kernel-based features using the simplex algorithm for face recognition [281], the extraction of facial subgraphs for head pose estimation [282], the analysis of Gabor kernels using univariate statistical techniques for discriminative region finding [283], the weighting of elastic graph nodes using quadratic optimization for authentication [284], the use of principal component analysis (PCA) to determine the importance of Gabor features [285], boosting Gabor features [286] and Gabor frequency/orientation selection using genetic algorithms [287].

A relevant work on Gabor Filters for face recognition that is closely related to the research presented here is by Wiskott and von der Malsburg [288]. Their work [289] [290] [291] [292], [288] proposes a framework for face recognition that is based on modeling human face images as labeled graph. Termed *Elastic Bunch Graph Matching* (EBGM), the technique has become a cornerstone in face recognition research. Each node of the graph is represented by a group of Gabor filters/wavelets (called "jets") which are used to model the intensity variations around their locations. The edges of the graph are used to model the relative location of the various jets. Since the jets represent the underlying image characteristics, it is desirable to place them on fiducial points on the face. This is achieved by *manually* marking the locations of the facial fiducial points using a small set of controlled graphs that represent "general face knowledge", which represents an average geometry for the human face. In our work, a genetic algorithm is used to obtain the spatial location of the fiducial points. Besides automating the process of locating these points, our work identifies spatial locations on the face image that are unique to every single person, rather than relying on an average geometry.

Closely following the work of Wiskott et al., Lyons et al. [293] proposed a technique that uses Gabor Filter coefficients extracted at 1) automatically located rectangular grid points or 2) manually selected image feature points. These coefficients are then used to bin face images based on sex, race and expression. The technique relies on a combined

Principal Component Analysis (PCA) dimensionality reduction and Linear Discriminant Analysis (LDA) classification over the extracted Gabor coefficients, to achieve a pooling of images. While the classification task is not related directly to *identifying* individuals from face images, this technique also demonstrates the ability of Gabor Filters to extract features that can encode subtle variations on facial images, providing a basis for face identification.

7.3.1.2 Gabor Filters

Mathematically, Gabor Filters can be defined as follows:

$$\Psi_{\omega,\theta}(x,y) = \frac{1}{2\pi\sigma_x\sigma_y} \cdot G_{\theta}(x,y) \cdot S_{\omega,\theta}(x,y) \quad (7.1)$$

$$G_{\theta}(x,y) = \exp \left\{ - \left(\frac{(x \cos \theta + y \sin \theta)^2}{2\sigma_x^2} + \frac{(-x \sin \theta + y \cos \theta)^2}{2\sigma_y^2} \right) \right\} \quad (7.2)$$

$$S_{\omega,\theta}(x,y) = \left[\exp \{ i(\omega x \cos \theta + \omega y \sin \theta) \} - \exp \left\{ - \frac{\omega^2 \sigma^2}{2} \right\} \right] \quad (7.3)$$

where,

- $G_{\theta}(x,y)$ represents a Gaussian Function.
- $S_{\omega,\theta}(x,y)$ represents a Sinusoid Function.
- (x,y) is the spatial location where the filter is centered with respect to the image axis.
- ω is the frequency parameter of a 2D Sinusoid.
- σ_{dir}^2 represents the variance of the Gaussian (and thus the filter) along the specified direction. *dir* can either be x or y . The variance controls the region around the center where the filter has influence.

From the definition of Gabor filters, as given in Equation 7.1, it is seen that the filters are generated by multiplying two components: a Gaussian Function $G_{\theta}(x,y)$ (Equation 7.2) and a Sinusoid $S_{\omega,\theta}(x,y)$ (Equation 7.3). The following discussions detail the two components of Equation 7.1.

7.3.1.3 Gaussian Function

The 2D Gaussian function defines the spatial spread of the Gabor filter. This spread is defined by the variance parameters of the Gaussian, along the x and y direction together with the orientation parameter θ . Figure 7.1(a) shows a 3D representation of the Gaussian mask generated with $\sigma_x = 10$ and $\sigma_y = 15$ and rotation angle $\theta = 0$. The image in Figure 7.1(b) shows the region of spatial influence of an elliptical mask on an image, where the variance in the x direction is larger than the variance in the y direction.

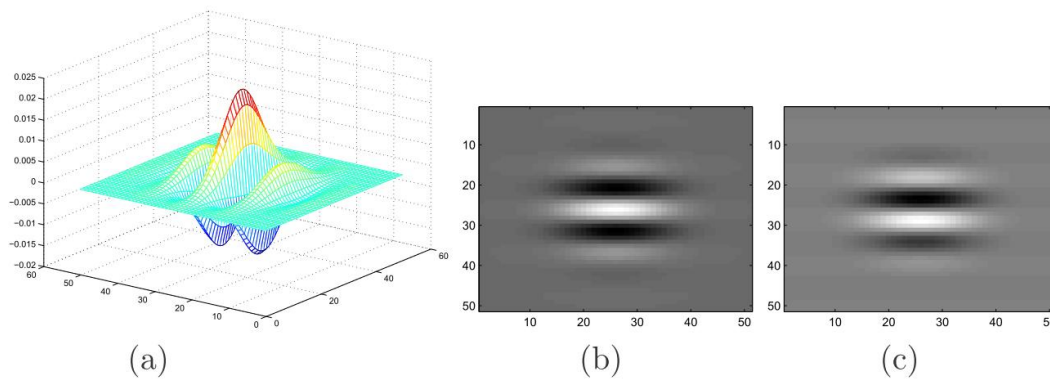


Figure 7.1: (a) 3D representation of a Gaussian mask; $\sigma_x = 10$, $\sigma_y = 15$ and $\theta = 0$
(b) Image of the Gaussian mask $\sigma_x = 10$, $\sigma_y = 15$ and $\theta = 0$

Typically the Gaussian filter has the same variance along both the x and y directions, that is $\sigma_x = \sigma_y = \sigma$. Under such conditions the rotation parameter θ does not play any role as the spread will be circular.

7.3.1.4 Sinusoid

The 2D complex Sinusoid defined by Equation 7.3 generates the two Sinusoidal components of the Gabor filters which (when applied to an image) extracts the local frequency content of the intensity variations in the signal. The complex Sinusoid has two components (the real and the imaginary parts) which are two 2D sinusoids that are phase shifted by $\frac{\pi}{2}$ radians. Figure 7.2(a) shows the 3D representation of a Sinusoidal signal (either real or imaginary) at $\omega = 0.554$ radians and $\theta = 0$ radians, while Figure 7.2(b) and 7.2(c) show an

image of the real and imaginary parts of the same complex Sinusoid, respectively. It can be seen that the two filters are similar, except for the π radian phase shift.

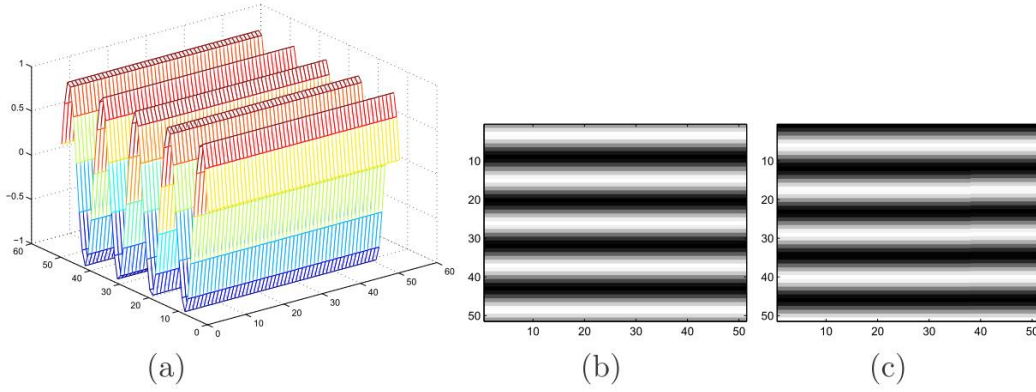


Figure 7.2: (a)3D representation of a Sinusoid $S_{\omega, \theta}$
 (b)Image representation of the real part of the complex Sinusoid $\Re \{S_{\omega, \theta}\}$
 (c)Image representation of the imaginary part of complex Sinusoid $\Im \{S_{\omega, \theta}\}$

Multiplying the Gaussian and the sinusoid generates the complex Gabor filter, as defined in Equation 7.1. If $\sigma_x = \sigma_y = \sigma$, then the real and imaginary parts of this complex filter can be described as follows.

$$\Re \{\Psi_{\omega, \theta}(x, y)\} = \frac{1}{2\pi\sigma^2} \cdot G_{\theta}(x, y) \cdot \Re \{S_{\omega, \theta}(x, y)\} \quad (7.4)$$

$$\Im \{\Psi_{\omega, \theta}(x, y)\} = \frac{1}{2\pi\sigma^2} \cdot G_{\theta}(x, y) \cdot \Im \{S_{\omega, \theta}(x, y)\} \quad (7.5)$$

Figure 7.3(a) shows the 3D representation of a Gabor filter (either real or imaginary) at $\omega = 0.554$ radians, $\theta = 0$ radians, and $\sigma = 10$ and Figure 7.3(b) and 7.3(c) show an image with the real and imaginary parts of the complex filter.

In order to extract a Gabor feature at a location (x, y) of an image I , the real and imaginary parts of the filter are applied separately to the same location in the image, and a magnitude is computed from the two results. Thus, the Gabor filter coefficient at a location

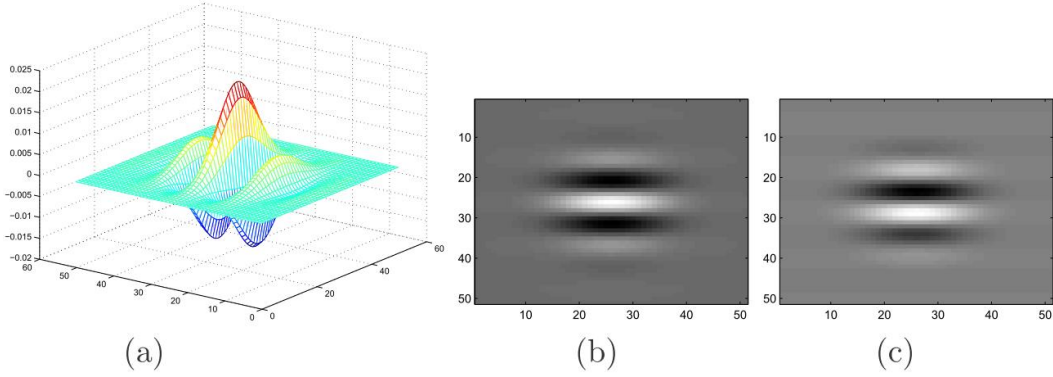


Figure 7.3: (a) 3D representation of a Gabor filter $\Psi_{\omega, \theta}$
 (b) Image representation of the real part of Gabor filter $\Re\{\Psi_{\omega, \theta}\}$
 (c) Image representation of the imaginary part of Gabor filter $\Im\{\Psi_{\omega, \theta}\}$

(x, y) in an image I with a Gabor filter $\Psi_{\omega, \theta}$ is given by

$$C_{\Psi}(x, y) = \sqrt{(I(x, y) * \Re\{\Psi_{\omega, \theta}(x, y)\})^2 + (I(x, y) * \Im\{\Psi_{\omega, \theta}(x, y)\})^2} \quad (7.6)$$

In our experiments, a *Gabor filter bank* was created by varying three parameters of $\Psi_{\omega, \theta}$: (1) the frequency parameter ω , (2) the orientation parameter θ , and (3) the variance parameter σ . We chose five values for each of these parameters thereby generating 125 different Gabor filters.

- $\omega = (2^{(-f+2)/2} \cdot \pi)$ where, $f = \{0, 1, 2, 3, 4\}$
- $\theta = (\frac{\pi}{2} \cdot \frac{1}{5} \cdot t)$ where, $t = \{0, 1.25, 2.5, 3.75, 5\}$
- $\sigma = \{5, 10, 15, 20, 25\}$

7.4 The Learning Algorithm

The proposed method uses the above described Gabor filters to find distinguishing features (and corresponding feature locations) within a face image. That is, for each person in the database, the algorithm finds a set of Gabor filters which, when applied at their corresponding (x, y) locations within the image will produce coefficients that are unique for that individual. This means that all of the 125 Gabor filters in the filter bank are applied at each and

every location of each of the individual's face images, and then tested for their ability to distinguish every individual. Given a 128×128 face image, there will be $128 \times 128 \times 125 \times n$ filter coefficients that will be generated per face image per person, where n is the number of characteristic features to be extracted for each person. This must be computed for every person in the training set, which further increases the search space. To search such a vast space of parameter values (the size of the Gaussian mask, the frequency of the complex sinusoid, the orientation of the entire Gabor filter, and the (x, y) location where the filter is placed) it is important that some scheme for effective search be incorporated into the system. To this end, we have chosen Genetic Algorithms to conduct the search. For each person in the training set, all of the face images that depict to that person are indexed as positives, while all of the other face images in the database are indexed as negatives. Dedicated Genetic Algorithm based search is conducted with these positive and negative images, with the aim of finding a set of Gabor filters and filter locations that distinguish all the positives from the negatives.

7.4.1 Genetic Algorithms

When the parameter space is vast (as it is in our case) a Genetic Algorithm (GA) searches for the optimum solution by randomly picking parameter sets and evolving newer ones from the best performers. This happens over many generations, hopefully resulting in the optimum set of parameters. To start the search, the GA generates a random set of *parents*. Each parent is characterized by the presence of a *chromosome*. The chromosome internally encodes all the parameters that are used by the parent to perform the intended operation. In our case, the intended operation is face recognition. The parent uses the parameters that are found in its chromosome to derive the Gabor features on the positive and negative images.

Based on the ability of these features to distinguish a face from all others in the database, the parent is ranked within its population. This rank is also referred to as the *fitness of the parent*. The ranking of all the parents, based on their fitness, marks the end of a generation, and a new generation needs to be created. New generations are formed based on three important aspects of GAs, *Retention*, *Cross Over* and *Mutation*. A portion

of the newer generation is derived from the older generation, using the above mentioned methods, and the rest of the new generation is created randomly, maintaining the same overall number of parents between generations. Once a new population has been formed, the process of ranking parents occurs (as explained earlier) and a new generation is born out of that ranking. This iterative process continues until the parents in a certain generation are fit enough to achieve the given task (with the desired amount of success) or until the desired number of generations have evolved.

7.4.1.1 Use of Genetic Algorithms in Face Recognition

GAs have been used in face recognition to search for optimal sets of features from a pool of potentially useful features that have been extracted from the face images. Liu et al. [294] used a GA along with Kernel Principal Component Analysis (KPCA) for face recognition. In their approach, KPCA was first used to extract facial image features. After feature extraction using the KPCA, GAs were employed to select the optimal feature subset for recognition - or more precisely the optimal non-linear components. Xu et al. [295] used GAs along with Independent Component Analysis to recognize faces. After obtaining all the independent components using the Fast ICA algorithm, a genetic algorithm was introduced to select optimal independent components.

Wong and Lam [296] proposed an approach for reliable face detection using genetic algorithms with eigenfaces. After histogram normalization of face images and computation of eigenfaces, the 'k' most significant eigenfaces were selected for the computation of the fitness function. The fitness function was based on the distance between the projection of a test image and that of the training-set face images. Since GAs are computationally intensive, the search space for possible face regions was limited to possible eye regions alone.

Karungaru et al. [297] performed face recognition using template matching. Template matching was performed using a genetic algorithm to automatically test several positions around the target, and to adjust the size of the template as the matching process

progressed. The template was a symmetrical T-shaped region between the eyes, which covered the eyes, nose and mouth.

Ozkan [298] used genetic algorithms for feature selection in face recognition. In this work, the Scale Invariant Feature Transform (SIFT) [299] was used to extract features. Since SIFT was originally designed for object recognition in general, genetic algorithms were used to identify SIFT features, which are more suitable to face recognition.

Huang and Weschler [300] developed an approach to identify eye location in face images using navigational routines, which were automated by learning and evolution using genetic algorithms. Specifically, eye localization was divided into two steps: (i) the derivation of the saliency attention map, and (ii) the possible classification of salient locations as eye regions. The saliency map was derived using a consensus between navigation routines that were encoded as finite state automata (FSA) exploring the facial landscape and evolved using genetic algorithms (GAs). The classification stage was concerned with the optimal selection of features and the derivation of decision trees for confirmation of eye classification using genetic algorithms.

Sun and Yin [301] applied genetic algorithms for feature selection in 3D face recognition. An individual face model was created from a generic model and two views of a face. Genetic algorithms were used to select optimal features from a feature space composed of geometrical structures, the labeled curvature types of each vertex in the individualized 3D model.

Sun et al. [302] approached the problem of gender classification using a genetic algorithm to select features. A genetic algorithm was used to select a subset of features from a low-dimensional representation, which was obtained by applying PCA and removing eigenvectors that did not seem to encode information about gender.

As is evident from these citations, many feature-based approaches towards face recognition use genetic algorithms for feature selection. However, these approaches employ a single feature space derived from a set of face images. We believe that it is more

effective to employ aimed at extracting person-specific features, and that an effective way to do this is by using genetic algorithms. As observed by [279], humans initially learn to recognize faces based on person-specific characteristic features. This suggests that better recognition performance might be achieved by representing each person's face in a person-specific feature space that is learned using GAs.

The following paragraphs describe how we employed GAs to solve the problem of finding person-specific Gabor features aimed at face recognition.

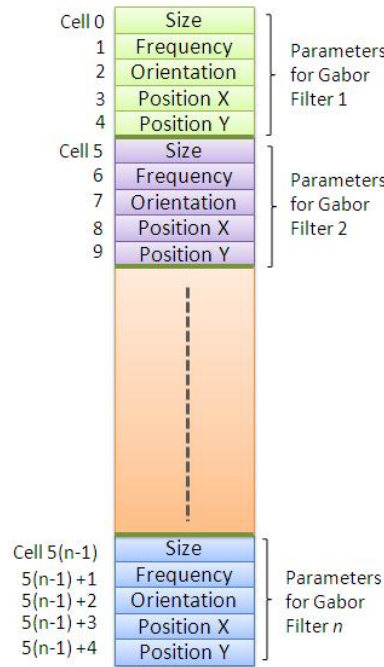


Figure 7.4: A typical chromosome used in the proposed method.

7.4.1.2 The Chromosome

Each parent per generation encodes the parameters of a set of Gabor filters in the form of a chromosome. In our implementation, each Gabor filter is represented by five parameters. If there are n Gabor filters, parameters for all of these filters are encoded into the chromosome in a serial manner, as shown in Figure 7.4. Thus the length of the chromosome is $5n$. The number of Gabor filters being used per face image determines the length of the chromosome. As shown in Figure 7.4, each parameter in the chromosome is encoded as a

gene. The boundaries of these genes defines the regions where the chromosome undergoes both the crossover and mutation. The genes can be considered as the primary element of the parent responsible in the evolution.

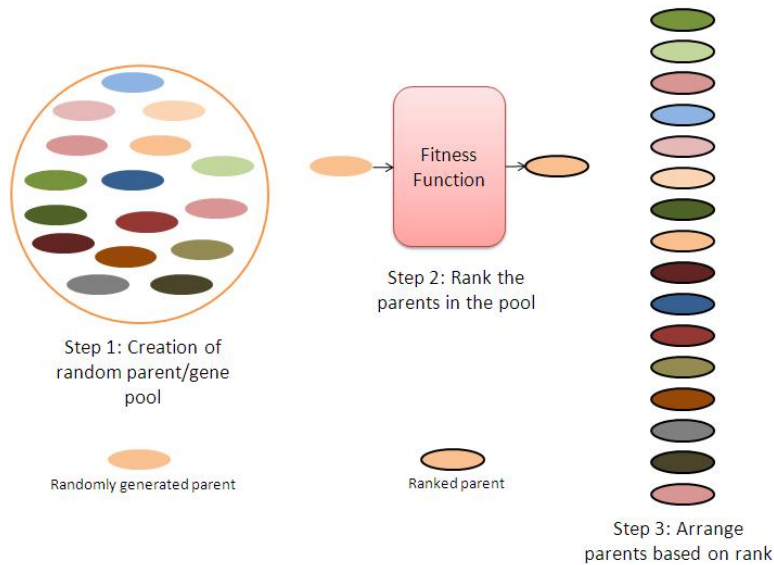


Figure 7.5: Stages in the creation of the first generation of parents

7.4.1.3 Creation of the first generation

Figure 7.5 depicts the first generation of parents, which are created randomly. Each parent's chromosome is filled randomly with parameter values where, each parameter value is within the allowed range for that parameter. Thus, in our experiment, each parent potentially has the parameters needed for it to perform face recognition using Gabor filters for feature extraction.

Once these parents are created, each parent in the gene pool is evaluated based on its capacity to perform face recognition. To this end, a fitness function is defined, which takes into account the ability of each parent to distinguish an individual from all others based on the most distinguishing features on the individual's face.

This fitness function also takes into account the similarity of the extracted features, and discourages the selection of features that are highly correlated with each other. This

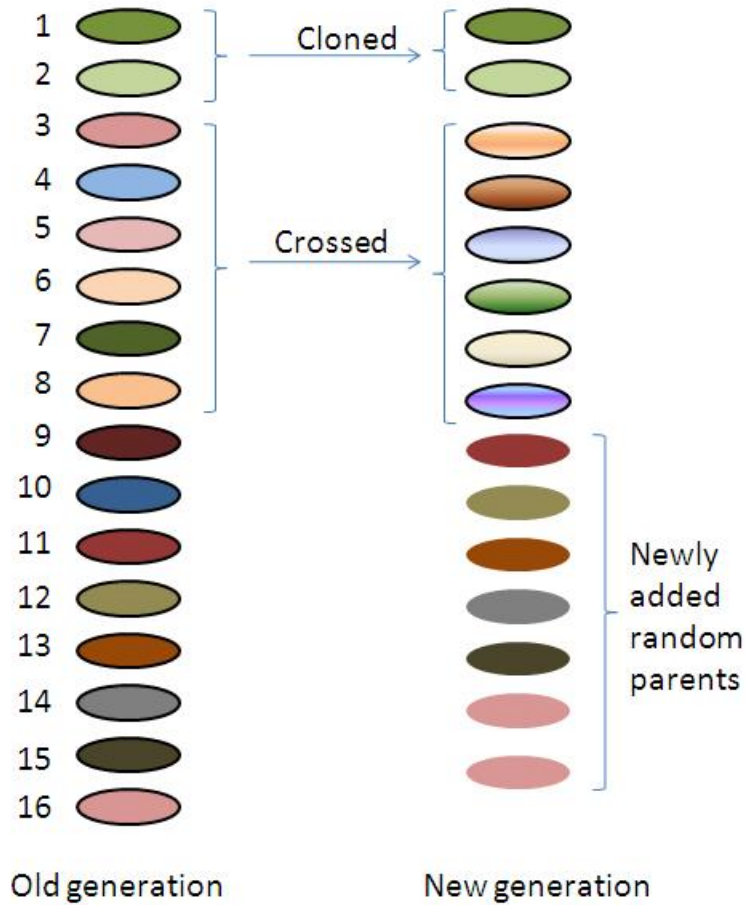


Figure 7.6: Deriving newer parents from the current generation

ensures that the face images will be searched for multiple distinguishing characteristics. Subsection 7.5.3.1 explains in detail the fitness function used in our experiments. The parents with the best fitness are ranked higher, and have the highest probability of being picked for using genetics the next generation. At the end of the rank ordering process, the parents are arranged in a descending order, based on their fitness. This rank ordering determines the probability of each parent being used to create the subsequent generation. If a parent has a higher fitness, it will have a higher probability of being cloned into the next generation, or of otherwise being involved in reproduction.

7.4.1.4 Creation of the newer generations

The newer generations are created from the older population using *clones*, *mutants*, and *crossovers* of the fittest parents. To better search for the optimal parameter set, new random parents are created every generation. This reduces the likelihood that the algorithm will get stuck in a local minimum in the search space.

Figure 7.6 shows crossover creates a newer generation, using the fittest parents from the older generation.

The number of offsprings created from mutation, cloning, and crossover are determined by parameters of the Genetic algorithm. The number of clones, mutants, and crossovers are controlled by the following parameters:

1. *Cloning Rate* This parameter controls the number of parents from the previous generation that will be retained without undergoing any changes in their genetic structure.
2. *Crossover Rate* This parameter controls the number of offsprings that will be born from crossing the parents from the previous generation.
3. *Mutation Rate* This parameter determines how many of the crossed offsprings will then be mutated.
4. *Cloning Distribution Variance* After determining the number of offsprings to be cloned, the index of the parents for cloning are chosen using a normal distribution random number generator, with the mean zero and variance equal to this parameter. Since the parents from the previous generation have been rank ordered in descending order of fitness, the zeroth parent will be the top performer (which coincides with the mean of the random number generator, and has the highest probability of getting picked).
5. *Crossover Distribution Variance* This parameter (which is similar to the Cloning Distribution Variance) is used to choose the index of the parents who will undergo

Crossover.

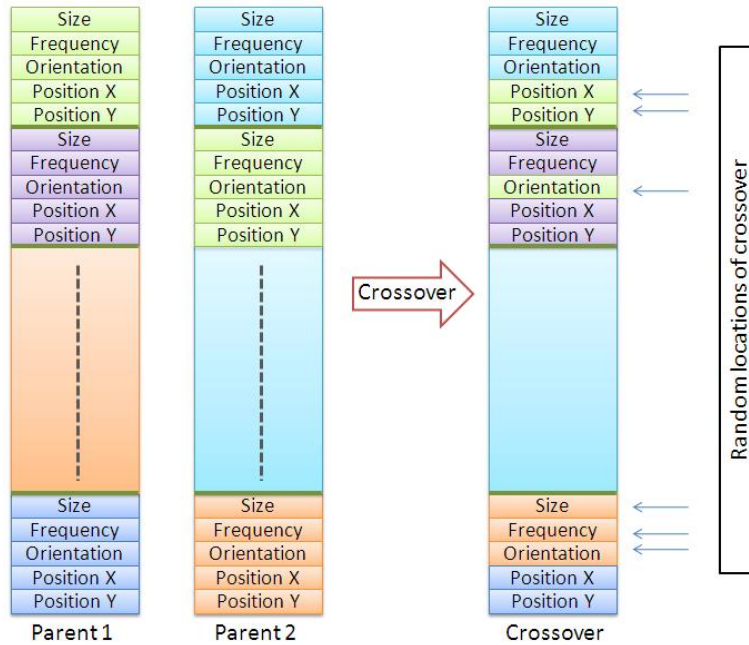


Figure 7.7: Typical crossing of two parents to create an offspring

7.4.1.5 Crossover

As discussed earlier, the parents for crossover are selected by a random number generator. Between these parents, the points of crossover are determined by choosing locations of crossover randomly. As seen in the Figure 7.7, these locations are arbitrary gene boundary locations and at these locations the gene content from the two parents gets mixed. The offspring thus created now contains parts of the genes coming from the contributing parents. The motivation for this step is the fact that, as more and more generations pass, the fittest parents undergoing crossover will already contain the better sets of parameters, and their crossing might bring together the better sets of parameter values from both the parents.

7.4.1.6 Mutation

In addition to the process of crossover at gene boundaries in the chromosome, the values of some parameters within the genes might be changed randomly. This is illustrated in the

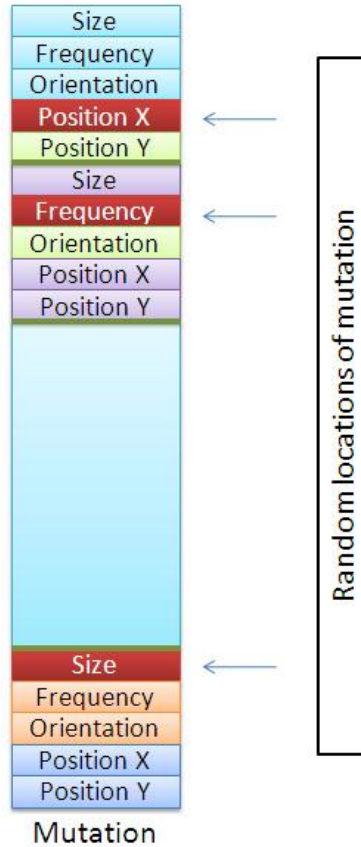


Figure 7.8: Mutation of a newly created offspring

Figure 7.8. Such mutations help in exploring the local parameter space more thoroughly. Mutations can be seen as small perturbations to the larger search that explores the vast parameter space, searching for the global minima.

7.5 Methodology

Most feature-based face recognition methods use feature detectors that are not tailored specifically for face recognition, and they make no attempt to selectively choose feature detectors based specifically on their usefulness for face recognition. The method described in this paper uses Gabor wavelets as feature detectors, but evaluates the usefulness of each particular feature detector (and a corresponding (x,y) location) for distinguishing between the faces within our face database. Given the very large number of possible Gabor feature detectors and locations, we use a Genetic Algorithm (GA) to explore the space of possibil-

ities, with a fitness function that propagates parents with a higher ability to distinguish between the faces in the database. By selecting the Gabor feature detectors and locations that are most useful for distinguishing each person from all of the other people in the database, we define a unique (i.e. person-specific) feature space for each person.

7.5.1 *The FacePix (30) Database*

All experiments were conducted with face images from the FacePix (30) database [303]. FacePix(30) was compiled to contain face images with pose and illumination angles annotated in 1 degree increments. Figure 7.9 shows the apparatus that is used for capturing the face images. A video camera and a spotlight are mounted on separate annular rings, which rotate independently around a subject seated in the center. Angle markings on the rings are captured simultaneously with the face image in a video sequence, from which the required frames are extracted.



Figure 7.9: The data capture setup for FacePix(30)

This database has face images of 30 people across a spectrum of pose and illumination

mination angles. For each person in the database, there are three sets of images. (1) The *pose angle set* contains face images of each person at pose angles from $+90$ to 90 (2) The *no-ambient-light set* contains frontal face images with a spotlight placed at angles ranging from $+90$ to -90 with no ambient light, and (3) The *ambient-light set* contains frontal face images with a spot light placed at angles placed at angles from $+90$ to -90 in the presence of ambient light. Thus, for each person, there are three face images available for every angle, over a range of 180 degrees. Figure 7.10 provides two examples extracted from the database, showing pose angles and illumination angles ranging from -90 to $+90$ in steps of 10. For earlier work using images from this database, please refer [304]. Work is currently in progress to make this database publicly available.



Figure 7.10: Sample face images with varying pose and illumination from the FacePix(30) database

We selected at random two images out of each set of three frontal (0) (Figure 7.11) images for training, and used the remaining image for testing. The genetic algorithms used the training images to find a set of Gabor feature detectors that were able to distinguish each persons face from all of the other people in the training set. These feature detectors were then used to recognize the test images.

In order to evaluate the performance of our system, we used the same set of training and testing images with face classification algorithm based on low-dimensional representation of face images extracted through Principal Component Analysis [305]. Specifically, the performance of the implementation of PCA-based face recognition followed by [306]



Figure 7.11: Sample frontal images of one person from the FacePix(30) Database

was used in our experiments.

7.5.2 *The Gabor Features*

Each Gabor feature corresponds to a particular Gabor wavelet (i.e. a particular spatial frequency, a particular orientation, and a particular Gaussian-defined spatial extent) applied to a particular (x, y) location within a normalized face image. (Given that 125 different Gabor filters were generated, by varying ω , σ and θ in 5 steps each, and given that each face image contained $128 \times 128 = 16,384$ pixels, there was a pool of $125 \times 16384 = 2,048,000$ potential Gabor features to choose from.) We used an N-dimensional vector to represent each person's face in the database, where N represents the predetermined number of Gabor features that the Genetic Algorithm selected from this pool. Figure 7.12 shows an example face image, marked with 5 locations where Gabor features will be extracted (i.e. $N = 5$). Given any normalized face image, real number Gabor features are extracted at these locations using Equation 7.6. This process can be envisioned as a projection of a 16,384-dimensional face image onto an N dimensional subspace, where each dimension is represented by a single Gabor feature detector.

Thus, the objective of the proposed methodology is to extract an N dimensional real-valued person-specific feature vector to characterize each person in the database. The N (x, y) locations (and the spatial frequency and spatial extent parameters of the N Gabor wavelets used at these locations) are chosen by a GA, with a fitness function that takes into

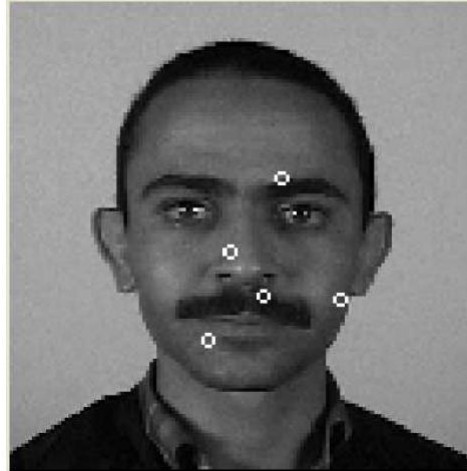


Figure 7.12: A face image marked with 5 locations where unique Gabor features were extracted

account the ability of each Gabor feature detector to distinguish one face from all the other faces in the database.

7.5.3 *The Genetic Algorithm*

Every GA is controlled in its progress through generations with a few control parameters such as,

- the number of generations of evolution (n_g)
- the number of parents per generation (n_p)
- the number of parents cloned per generation (n_c)
- the number of parents generated through cross over (n_{co})
- the number of mutations in every generation (n_m)

In our experiments, the GA used the following empirically-chosen GA parameters:

$n_g = 50$, $n_p = 100$, $n_c = 6$, $n_{co} = 35$ and $n_m = 5$.

7.5.3.1 The Fitness Function

The fitness function of a genetic algorithm determines the nature of the search conducted over the parameter space. For face recognition applications, the fitness function is the capacity of a parent to classify the individuals accurately. In our proposed method, the fitness function needs to take both the Gabor features and the corresponding feature locations into consideration when evaluating face classification. We define here a fitness function that has two components to it. One determines the capacity of the parent to isolate an individual's face image from the others in the database, and the other evaluates whether the feature is redundant with other extracted features (i.e. whether a feature detector produces coefficients that are highly correlated with the coefficients produced by another feature detector.) Thus the fitness F can be defined as

$$F = w_D D - w_C C \quad (7.7)$$

where D is the distance measure weighted by w_D , and C represents the correlation measure which measure the similarity between the coefficients that have been extracted. The correlation measure C is weighted by the factor w_C .

If a parent extracts features from a face image that distinguish one individual from all the others very well (compared to the other parents within the same generation) then the distance measure D will be the largest for that parent, making its fitness F large. If the correlation between the extracted features is small, C will be small, which also makes the fitness F large. Thus, the correlation measure serves as a *penalty* for extracting the same feature from the face image multiple times, even though that particular feature might be the best distinguishing feature on that face.

The correlation between coefficients was used instead of spatial separation to counter the problem of similar features being extracted, because the Gabor filters might not be able to represent the underlying image characteristic completely. If there are some large image

features on the face (such as beard) that require multiple Gabor features within a certain spatial locality. Setting a hard lower limit on this spatial separation might lead to insufficient representation of that large image feature, in terms of the Gabor filters.

Consider a parent searching for a unique set of M Gabor filters to distinguish one individual's face from all other faces. Let this set of filters be referred to as S . Thus, $S = \{G_1, G_2, \dots, G_M\}$ where, G_m represents the m^{th} Gabor filter.

If the set all individuals in the database is referred to as $I = \{i_1, i_2, \dots, i_j\}$ with J number of individuals, then for every individual i in I a set S_i has to be extracted. To achieve this, all the images in the database depicting individual i are marked as positives, and the ones not depicting that individual are marked as negatives. Let the set of positive images be referred to as P_i (with L number of images) and the set of negatives be referred to as N (with K number of images). Thus, $S_i = \{G_{1i}, G_{2i}, \dots, G_{mi}\}$, $P_i = \{p_{1i}, p_{2i}, \dots, p_{li}\}$ and $N_i = \{n_{1i}, n_{2i}, \dots, n_{ki}\}$ are the sets of Gabor filters, positive images and negatives images set respectively for the individual i .

- **The Distance Measure D**

A parent trying to recognize an individual i with a Gabor filter set S_i can be thought of as a transformation that projects all of the face images from the image space to a M -dimensional space, where the dimensions are defined by the M Gabor filters in the set S_i . Thus, all of the images in the two sets P_i and N_i can be considered as points on this M -dimensional space. Since the goal of the genetic algorithm is to find the set S_i which best distinguishes the individual i from others, in our method we search for the M dimensional space (defined by a parent) that best separates the points formed by the sets P_i and N_i . Figure 7.13 is an illustration of hypothetical set of face images projected on a 2 dimensional space defined by a set of 2 Gabor filters $S_i = \{G_0, G_1\}$. As shown in the figure, the measure D is the minimum of all the Euclidian distances between every positive and negative points.

Thus, D can be defined as follow:

$$D = \min_{\forall i,k} [\delta_M(\phi_M(p_{li}), \phi_M(n_{ki}))] \quad (7.8)$$

where,

$\delta_M(A,B) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_m - b_m)^2}$ is the M -dimensional Euclidian distance between A and B . a_x and b_x corresponds the x^{th} -coordinate of A and B respectively

$\phi_M(X)$ is the transformation function that projects image X from the image space to the M -dimensional space defined by the set of Gabor filters.

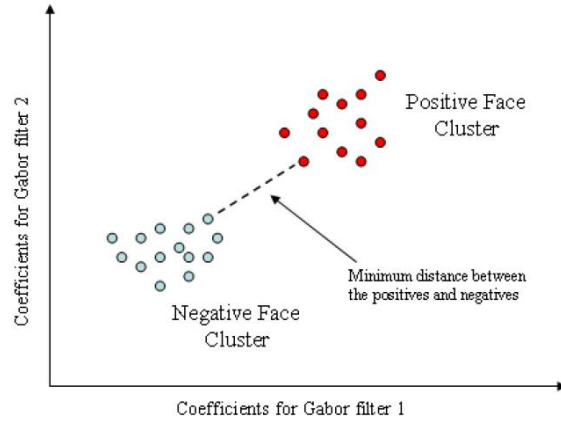


Figure 7.13: Distance Measure D for the fitness function

- **The Correlation Measure C**

In the proposed method, in addition to having every parent selecting the Gabor filter set S_i that can best distinguish the individual i from all the others in the database, it is necessary to ensure that this set of Gabor filters does not include filters that extract identical image features. If there were no such constraint, the algorithm might find one very distinguishing image feature on the face image and, over generations of evolution, all of its Gabor filters might converge to this one image feature. To avoid this, the correlation measure C determines the correlation between the image features extracted at all the locations pointed to by the chromosome. To test for correlations

between the Gabor features at the different spatial locations, we use the entire set of 125 Gabor filters to thoroughly characterize the textural context at these locations.

Assuming that there are M Gabor features that we are looking for on the face image of individual i , let $(x_m, y_m), m = 1, 2, \dots, M$ be the M points that have been selected genetically in the chromosome. To find the correlations of the image features extracted at each of these points, the N Gabor filters $G_i, i = 1, 2, \dots, N$ are used to characterize each of the points. Let the coefficients of such a characterization be represented by a matrix A . Thus, matrix A is $M \times N$ in dimension, where the rows correspond to the M locations and $N = 125$ refers to the Gabor filter coefficients. Thus,

$$A = \begin{bmatrix} g(1,1) & g(1,2) & \cdots & g(1,N) \\ g(2,1) & g(2,2) & \cdots & g(2,N) \\ \vdots & \vdots & \vdots & \vdots \\ g(m,1) & g(m,2) & \cdots & g(m,N) \end{bmatrix} \quad (7.9)$$

where, $g_{(m,n)}$ is the coefficient obtained by applying the n^{th} Gabor filter to the image at the point (x_m, y_m) .

The Correlation measure can now be defined in terms of matrix A as follows

$$C = \log(\det(diag(B))) - \log(\det(B)) \quad (7.10)$$

where, $diag(B)$ returns the diagonal matrix corresponding to B , and B is the covariance matrix defined by $B = \frac{1}{N-1}(AA^T)$.

Examining the Equation 7.10, it can be seen that the first log term gets closer to the second log term when the off diagonal elements of B reduces. The diagonal elements of the matrix B corresponds to the variance of the M image locations, whereas the off diagonal elements correspond to the covariance between pairs of locations. Thus, as the covariance between the image points decreases, the value of the overall correlation parameter decreases.

- **Normalization of D and C**

In order to have an equal representation of both the Distance measure D and the Correlation term C in the fitness function, it is necessary to normalize the range of values that they can take. For each generation, before the fitness values are used to rank the parents, parameters D and C are normalized to range between 0 and 1.

$$D_{norm} = \frac{D - D_{Min}}{D_{Max} - D_{Min}} \quad (7.11)$$

$$C_{norm} = \frac{C - C_{Min}}{C_{Max} - C_{Min}} \quad (7.12)$$

where, the Max represents the maximum value of D or C in a single generation across all the parents and Min refers to the minimum value.

- **Weighting factors w_D and w_C**

The influence of the two components of the fitness function are controlled by the weighting factors w_D and w_C . We used the relation $w_C = 1 - w_D$ to control the two parameters simultaneously. With this relationship, a value of $w_D \approx 1$ will subdue the effect of the Correlation measure, causing the genetic algorithm to choose the Gabor filters on the most prominent image feature alone. On the other hand, $w_D \approx 0$ will subdue the Distance measure, deviating the genetic algorithm from the main goal of face recognition. Thus an optimal value for the weight w_D has to be estimated empirically, to suit the face image database in question.

7.6 Results

To evaluate the relative importance of the two terms (D and C) in the fitness function, we ran the proposed algorithm on the training set several times with 5 feature detectors per chromosome, while changing the weighting factors in the fitness function for each run, setting w_D to 0, .25, .50, .75, and 1.00, and computing $w_C = (1 - w_D)$. Figure 7.14 shows the recognition rate achieved in each case.

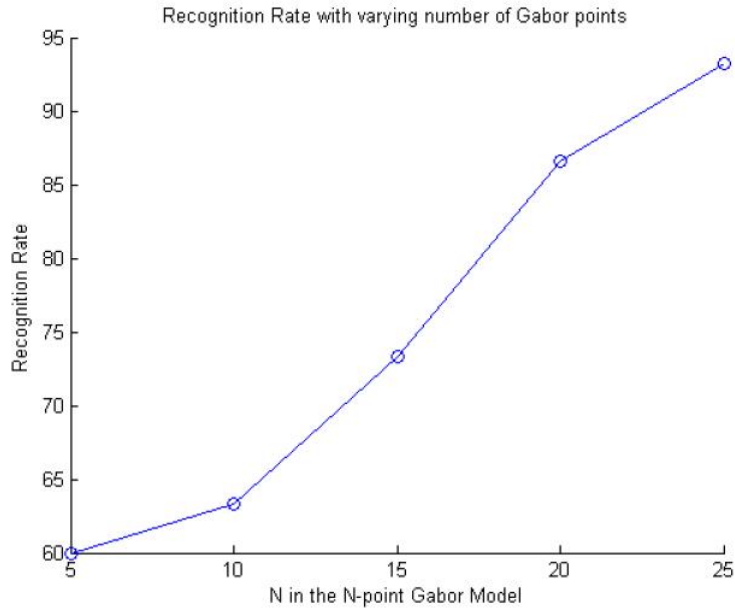


Figure 7.14: The recognition rate versus the number Gabor feature detectors

We then ran the proposed algorithm on the training set 5 times, while changing the number of Gabor feature detectors per parent chromosome for each run to 5, 10, 15, 20, and 25. In all the trials, $w_D=0.5$. Figure 7.15 shows the recognition rate achieved in each case.

7.6.1 Discussion of Results

Figure 7.14 shows that the recognition rate of the proposed algorithm when trained with 5, 10, 15, 20, and 25 Gabor feature detectors increases monotonically, as the number of Gabor feature detectors (N) is increased. This can be attributed to the fact that increasing the number of Gabor features essentially increases the number of dimensions for the Gabor feature detector space, allowing for greater spacing between the positive and the negative clusters.

Figure 7.15 shows that for $N = 5$ the recognition rate was optimal when the distance measure D and the correlation measure C were weighted equally, in computing the fitness function F . The dip in the recognition rate for $w_D = 0.75$ and $w_D = 1.0$ indicates the significance of using the correlation factor C in the fitness function. The penalty introduced

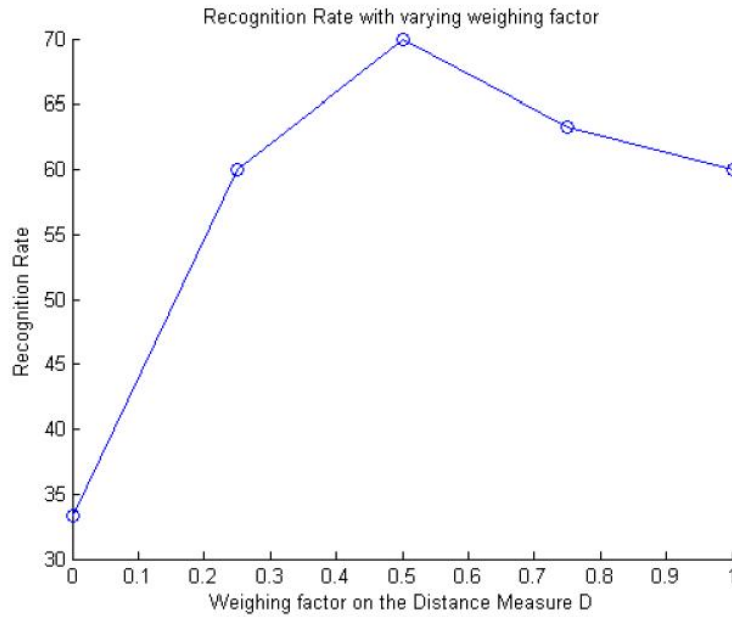


Figure 7.15: Recognition rate with varying w_D

by C ensures that the GA searches for Gabor features with different textural patterns. If no such penalty were be imposed, the GA might select Gabor features that are clustered on one salient facial feature, such as a mole.

The best recognition results for the proposed algorithm (93.3%) were obtained with 25 Gabor feature detectors. The best recognition performance for the PCA algorithm was reached at about 15 components, and flattened out beyond that point, providing a recognition rate for the same set of faces that was less than 83.3%. This indicates that, for the face images used in this experiment (which included substantial illumination variations) the proposed method performed substantially better than the PCA algorithm.

7.6.2 Person-specific feature extraction

When the FacePix(30) face database was built, all but one person were captured without eyeglasses or a hat. Figures 7.16(a) and 7.16(b) show the results of extracting 10 and 20 distinguishing features from that person's face images. The important things to note about these results are:

1. At least half of the extracted Gabor features (8 of the 10) and (10 of the 20) are located on (or near) the eyeglasses.
2. As the number of Gabor features was increased from 10 to 20, more Gabor features are seen toward the boundaries of the images. This is due to the fact that the genetic algorithm chooses Gabor feature locations based on a Gaussian probability distribution that is centered over the image, and decreases toward the boundaries of the images.

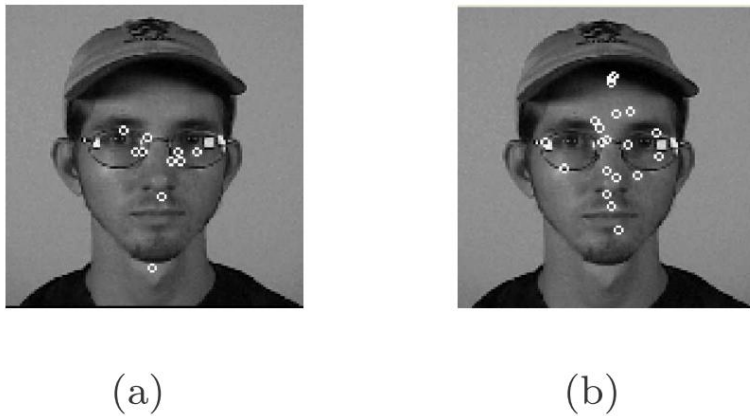


Figure 7.16: 10 and 20 person-specific features extracted for a particular individual in the database

These results suggest that person-specific feature extraction might be useful for face recognition in small face databases, such as those typical of a social interaction assistance device for people who are blind.

7.7 Face Recognition in Social Situational Awareness Context

As mentioned earlier, the proposed person-specific approach to evolutionary feature selection in face images is well-suited for applications such as those that enhance social interaction for people who are blind, because people do not generally disguise their appearance in normal social situations, and even when some significant change occurs (such as a man shaving off his beard) the system can continue to evolve as it captures new images with each encounter.

In summary, while there have been many different feature-based approaches to face recognition over the last two decades of research, we have proposed a novel methodology based on the discovery and extraction of person-specific characteristic features to improve face recognition performance for small face databases. This approach is aimed at facilitating social interaction in casual settings. The use of Gabor features, in tandem with a genetic algorithm to discover characteristic person-specific features has been inspired by the human visual system and is based on knowledge that has been developed about the process by which humans recognize faces. We believe that more needs to be learnt about human face recognition, and that as more is learnt, the knowledge can be put to use to develop more robust face recognition algorithms.

In the following three chapters, we discuss the problem of sensing and conveying Proxemics (interpersonal spaces), which is an important component of social interactions and social situational awareness. These chapters, together, provide sensing and delivering of the proxemics information.

Chapter 8

ENRICHING GROUP INTERACTIONS: SENSING PROXEMICS AND THE SOCIAL SPACE

An important aspect of social interactions is the relative spatial location of the interacting partners. Definition of what a culture is correlates very closely to the accepted norms of interpersonal distance and the behavioral attributes displayed by the interacting partners. In behavioral psychology, influences of interpersonal distances on social interactions between people have been studied for over four decades. The term proxemics, coined by Edward T. Hall, describes influence of interpersonal distances in animal and man [307]. The following list describes the American proxemic distances; note that such distances vary with culture and environment. A part of learning that visually able individuals go through is to understand what these norms are very quickly by observing interactions between other individuals who are familiar with the culture. People who are visually impaired might find the lack of access to this important information inconvenient. An important aspect of the technologies proposed in the next three chapters is to provide information to the users a means of accessing the social scene.

1. Intimate Distance (Close Phase): 0-6 inches
2. Intimate Distance (Far Phase): 6-18 inches
3. Personal Distance (Close Phase): 1.5-2.5 feet
4. Personal Distance (Far Phase): 2.5-4 feet
5. Social Distance (Close Phase): 4-7 feet
6. Social Distance (Far Phase): 7-12 feet
7. Public Distance (Close Phase): 12-25 feet
8. Public Distance (Far Phase): 25 feet or more

In Section 1.2.1.2 of Chapter 1, we described the importance of group interactions and a high level view of the dynamics in group interactions. An important aspect of this group interaction is communicating the interpersonal distances and dynamics of individual's movements within the interpersonal space. In this chapter we describe a computer vision approach towards localizing individuals in the interpersonal space of the user of the social situational awareness device.

From the sensing perspective we extract the social scene structure, through the use of face detection/person detection and tracking, towards determining the number of people in the user's visual field, where people are located relative to the user, coarse information related to gaze direction (pose estimation algorithms could be used to extract finer estimates of pose), and the approximate distance of the person from the user based on the size of the face image. The camera is mounted on a pair of glasses that allow the camera to be directed by the user in which ever direction possible. The camera also tracks individuals as they move in the space in front of the users. Through the use of detected faces and a tracking algorithm it is possible so sense all the individuals within the interaction space and their relative location and movement in front of the user. This chapter describes the challenges involved in detecting faces within a scene and determining their exact size in order to determine the distance to the user. The following chapter describes tracking individuals within the social space, especially when they are moving in front of the users.

8.1 Accurate Face Detection

Face detection has become an important first step towards solving plethora of other computer vision problems like face recognition, face tracking, pose estimation, intent monitoring and other face related processing. Over the years many researchers have come up with algorithms, that have over time, become very effective in detecting faces in complex backgrounds. Currently, the most popular face detection algorithm is the Viola-Jones [308] face detection algorithm whose popularity is boosted of by its availability in the open source computer vision library, OpenCV. Other popular face detection algorithms are identified in [309] and [126].

Most face detection algorithms learn faces by modeling the intensity distributions in upright face images. These algorithms tend to respond to face-like intensity distributions in image regions that do not depict any face as they are not contextually aware of the presence or absence of a human face. These spurious responses make the results unsuitable for further processing that requires accurate face images as inputs, such as the ones mentioned above. Figure 8.1 shows an example where a face detection algorithm detects two faces - one true and the other false.



Figure 8.1: An example false face detection.

The problem of false face detection has motivated some researchers to develop heuristic approaches aimed for validating the face detection results. Most of these heuristics integrate primitive context into the problem by searching for skin tone in the output subimages. However, this simple approach often fails to distinguish faces from non-faces, because face detectors often fail to center the cropping box precisely around the detected face. This produces a significant patch of skin colored pixels, but only a partial face. This centering problem can be dealt with by extracting the skin colored regions and comparing their shape to an ellipse. While such heuristics, are simple, and somewhat effective, their validation is not reliable enough to meet the needs of higher level face processing tasks.

Further, they do not provide a confidence metric for their validation.

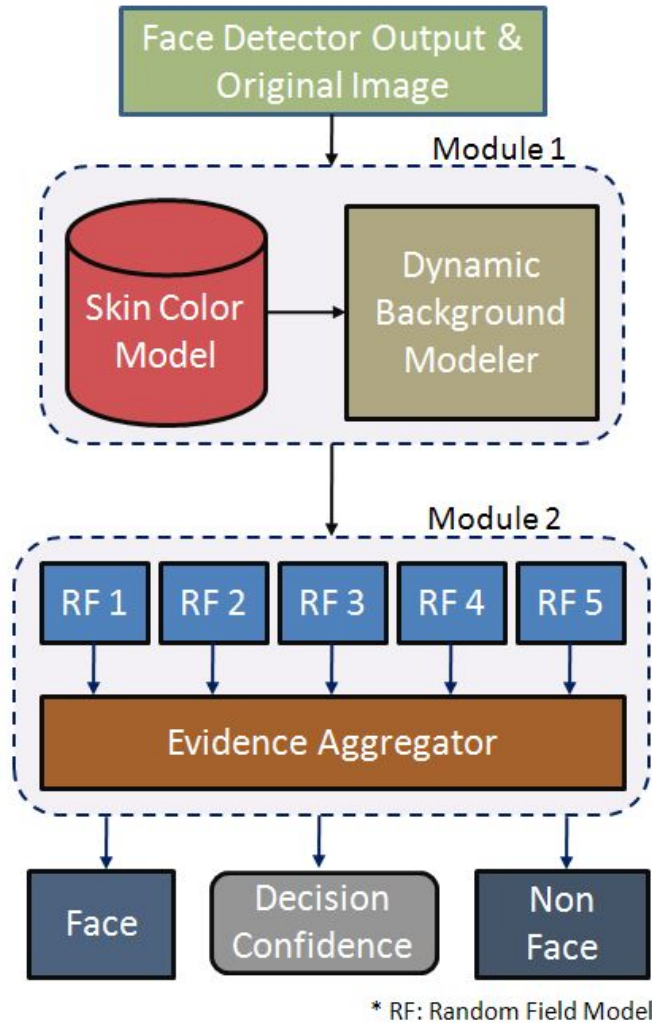


Figure 8.2: Block diagram.

This chapter treats the problem of face detection validation in a systematic manner, and proposes a learning framework that incorporates both contextual and structural knowledge of human faces. A face validation filter is designed by combining two statistical modelers, 1) a human skin-tone detector with a dynamic background modeler (Module 1), and 2) an evidence-aggregating human face silhouette random field modeler (Module 2), which provides a confidence metric on its validation task. The block diagram in Figure 8.2 shows the functional flow of data through the two modules in the proposed framework. The details of the statistical models and their learning will be presented later in the paper, which

is organized as follows. Section 2 reviews some of the earlier research. Section 3 introduces the proposed framework, with details on the learning process. Section 4 discusses the experiments carried out to test the proposed framework. Section 5 presents the results while Section 6 discusses them. Section 7 concludes the paper and discusses future work.

8.2 Related Work in Accurate Face Detection

As mentioned earlier, the problem of face detection validation has not been treated methodically before, though the problem has been handled by many as an integral component of face detection algorithms. All the past work in this area can be broadly characterized into two groups: a) Low level image feature models mostly based on skin color such as [310], [311] and [312], and b) High level facial feature models such as [313], [314] and [315].

The low level skin color based approaches try to reduce computational complexity by first identifying skin color in images so that search can be reduced. Most of the times, simple geometrical properties of the retained skin regions are used to determine if the region is a face. Such simplification of faces into trivial geometrical structures results in false detections. The facial feature based methods achieve face detection by individually identifying the integral components of a face image such as eyes, nose, etc. Though these schemes could be robust, the associated computational load is high. Interested readers could find more related references in [126] and [309]. The framework proposed in this paper uses statistically learnt knowledge about human faces to overcome computational complexity thereby augmenting face validation to existing face detection algorithms seamlessly.

8.3 Proposed Framework

As shown in Figure 8.2, the framework essentially has two statistically learnt models, Module 1 and Module 2, that are cascaded to form the face detection validation filter. The output from a face detector is sent to Module 1, which distinguishes the skin pixels in the face region from the background pixels, thereby constructing a skin region mask. This skin region mask then becomes the input to Module 2, which is essentially an aggregate of random field models learnt from manually labeled (*true*) face detection outputs. The results of each random field model within the aggregate are then combined, using rules of Dempster-Shafer

Theory of Evidence [316]. This combining of evidence provides a metric for the belief (i.e. confidence) of the system in its final validation. The two modules are detailed in the following subsections.

8.3.1 Module 1: Human Skin Tone Detector with Dynamic Background Modeler

Most of the skin tone detectors used for human skin color classification use prior knowledge, which is provided in the form of a parametric or non-parametric model of skin samples that are extracted from images - either manually, or through a semiautomated process. In this paper we employ such an a priori model, in combination with a dynamic background modeler, so that the skin vs. non-skin boundary is accurately determined. Accurate skin region extraction is essential for Module 2, as it validates images based on their structural properties. The two functional components of Module 1 are:

8.3.1.1 *a-priori* Bi-modal Gaussian Mixture Model for Human Skin Classification

A normalized RGB color space has been a popular choice among researchers for parametric modeling of human skin color. The normalized RGB (typically represented as nRGB) of a pixel X with X_r, X_g, X_b as its red, green and blue components respectively, is defined as:

$$X_{i|i \in \{r,g,b\}}^{nRGB} = \frac{X_i}{\left(\sum_{\forall i \in \{r,g,b\}} X_i \right)} \quad (8.1)$$

Normalized RGB space has the advantage that only two of the three components, nR, nG or nB, is required at any one time to describe the color. The third component can be derived from the other two as:

$$X_{i|i \in \{nR,nG,nB\}}^{nRGB} = 1 - \left(\sum_{\forall k \in \{nR,nG,nB\}, k \neq i} X_k \right) \quad (8.2)$$

In our experiments, we found that skin pixels form a tight cluster when projected on nG and nB space as shown in the Figure 8.3. The study was based on a skin pixel database, consisting of nearly 150,000 samples, built by randomly sampling skin regions from 1040 face images collected on the web as well as from FERET face database [317]. Further analysis also showed that the cluster formed on the 2D nG-nB space had two prominent density peaks which motivated the modeling of skin pixels with a Bi-modal Gaussian

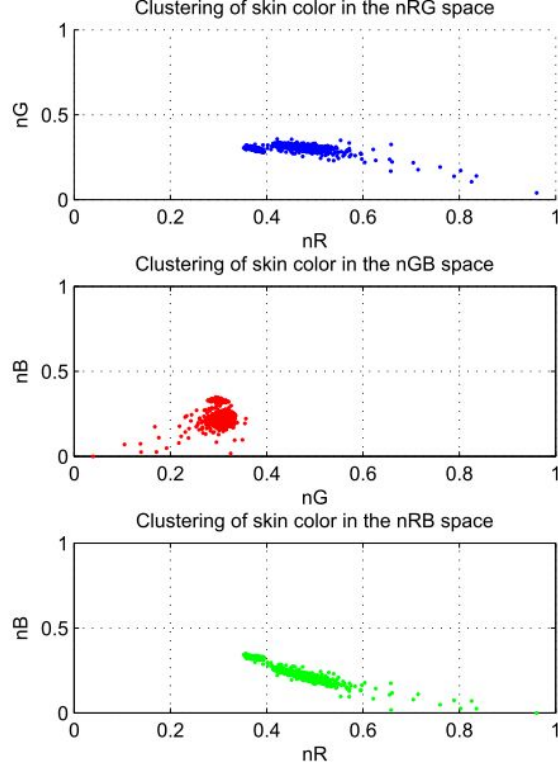


Figure 8.3: Skin pixels in nRGB space.

mixture model learnt using Expectation Maximization (EM) with a k -means initialization algorithm [318]. The Bi-modal Gaussian mixture model is represented as.

$$\begin{aligned}
 f_{X|X=[nG,nB]}^{skin}(x) &= w_1 f_{Y_1}(x; \Theta_1 = [\mu_1, \Sigma_1]) + \\
 &w_2 f_{Y_2}(x; \Theta_2 = [\mu_2, \Sigma_2])
 \end{aligned} \tag{8.3}$$

8.3.1.2 Dynamically Learnt Multi-modal Gaussian Model for Background Pixel

Classification

As mentioned earlier, classification of regions into face or non-face requires accurate skin vs. non-skin classification. In order to achieve this, we learn the background color surrounding each face detector output dynamically. To this end we extract an extra region of the original image around the face detector's output, as shown in Figure 8.4. Since the size of the face detector output varies from image to image, it is necessary to normalize the size. This is done by downsampling the size of the original image to produce a face detector out-

put region containing 90x90 pixels. The extra region pixels surrounding the face are then extracted from the 100x100 region around this 90x90 normalized face region.

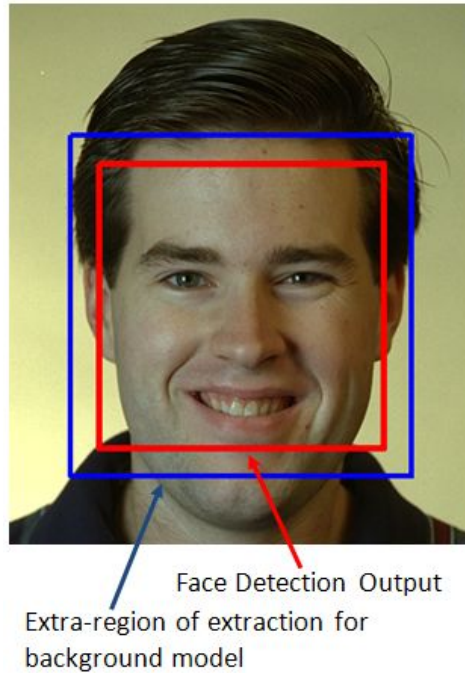


Figure 8.4: Extra region for background modeling.

Once the outer pixels are extracted, a Multi-modal Gaussian Mixture is trained using EM with k -means initialization, similar to the earlier case with skin pixel model. The resultant model can be represented as.

$$f_{X|X=[R,G,B]}^{non-skin}(x) = \sum_{i=1}^m w(i) f_{Y_i}(x; \Theta_i = [\mu_i, \Sigma_i]) \quad (8.4)$$

where, m is the number of mixtures in the model. We found empirically that a value of $m = 2$ or $m = 3$ modeled the backgrounds with sufficient accuracy.

8.3.1.3 Skin and Background Classification using the learnt Multi-modal Gaussian Models

The skin and non-skin models, $f_{X|X=[nG,nB]}^{skin}(x)$ and $f_{X|X=[R,G,B]}^{non-skin}(x)$ respectively, are used for classifying every pixel in the scaled face image obtained as explained in the Section 8.3.1.2. Example skin-masks are shown in Figure 8.5. This example shows two sets of images - one corresponding to a *true* face detection result, and another *false* face detection result.

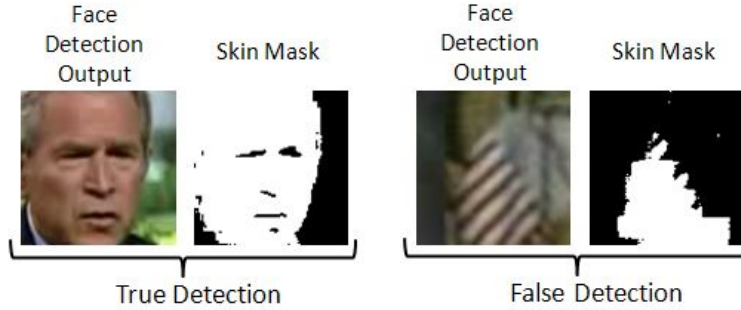


Figure 8.5: Example of *true* and *false* face detection.

The structural analysis through Random Field models explained in the next section will describe the design concepts that will help distinguish between *true* and *false* face detections shown in Figure 8.5.

8.3.2 Module 2: Evidence-Aggregating Human Face Silhouette Random Field Modeler

In order to validate the skin region extracted as explained in Section 8.3.1, we build statistical models from examples of faces. We developed statistical learners inspired by Markov Random Fields (MRF) to capture the variations possible in *true* skin masks (face silhouette). The following subsections describes MRF models and the variant we created for our experiments.

8.3.2.1 Random Field (RF) Models

In this work, we used a minor variant of MRFs to learn the structure of a *true* face skin mask. MRFs encompass a class of probabilistic image analysis techniques that rely on modeling the intensity variations and interactions among the image pixels. MRFs have been widely used in low level image processing including, image reconstruction, texture classification and image segmentation [319].

In an MRF, the sites in a set, \mathcal{S} , are related to one another via a neighborhood system, which is defined as $\mathcal{N} = \{\mathcal{N}_i, i \in \mathcal{S}\}$, where \mathcal{N}_i is the set of sites neighboring i , $i \notin \mathcal{N}_i$ and $i \in \mathcal{N}_j \iff j \in \mathcal{N}_i$.

A random field X said to be an MRF on \mathcal{S} with respect to a neighborhood system

\mathcal{N} , if and only if,

$$P(\mathbf{x}) > 0, \forall \mathbf{x} \in \mathcal{X} \quad (8.5)$$

$$P(x_i | x_{\mathcal{S}-\{i\}}) = P(x_i | x_{\mathcal{N}_i}) \quad (8.6)$$

where, $P(x_i | x_{\mathcal{S}-\{i\}})$ represents a Local Conditional Probability Density function defined over the neighborhood \mathcal{N} . The variant of MRF that we created for our experiments relaxed the constraints imposed by MRFs on \mathcal{N} . Typically, MRFs requires that sites in set \mathcal{S} be contiguous neighbors. The relaxation in our case allows for distant sites to be grouped into the same model.

We empirically found out that modeling the skin-region validation problem into one single RF gave poor results. We devised 5 unique RF models with an Dempster-Shafer Evidence aggregating framework that could not only validate the face detection outputs, but also provide a metric of confidence. Thus, Equation 8.6 could be alternatively seen as a set $P(\mathbf{x}) = \{P^1(\mathbf{x}), \dots, P^5(\mathbf{x})\}$, each having their own neighborhood system $\mathcal{N}^k = \{\mathcal{N}^1, \mathcal{N}^2, \dots, \mathcal{N}^5\}$, such that

$$P^k(x_i | x_{\mathcal{S}-\{i\}}) = P(x_i | x_{\mathcal{N}_i^k}) \quad (8.7)$$

8.3.2.2 Pre-processing

As described earlier, each face detector output is normalized and expanded to produce a 100x100 pixel image, from which a binary skin mask is generated. A morphological opening and closing operation is then performed on the skin mask (to eliminate isolated skin pixels), and the mask is then partitioned into one hundred 10x10 blocks, as shown in Figure 8.6. The number of mask pixels (which represent skin pixels) are counted in each block, and a 10x10 matrix is constructed, where each element of this matrix could contain a number between 0 and 100. This 10x10 matrix is then used as the basis for determining whether the face detector output is indeed a face.

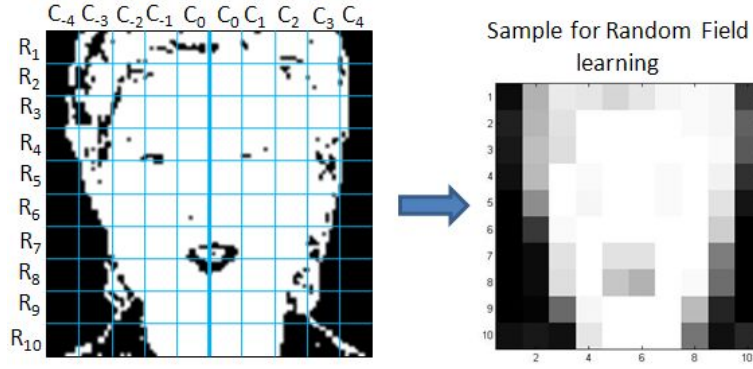


Figure 8.6: Pre-processing.

8.3.2.3 The Neighborhood System

The determination of whether the face detector output is actually a face is based on heuristics that are derived from anthropological human face models [320] and through our own statistical analysis. These include:

1. Human faces are horizontally symmetrical (i.e. along any row of blocks R_i) about a central vertical line joining the nose bridge, the tip of the nose and the chin cleft, as shown in Figure 8.6. In particular, our analysis of a large set of frontal face images showed that the counts of skin pixels in the 10 blocks that form each row in Figure 8.6 were roughly symmetrical across this central line.
2. The variations along the verticals (C_i 's) are negligible enough that in building a Local Conditional Probability Density function, each R_i can be considered independent of the other. That is, for example, modeling variations of C_0 w.r.t C_1 on R_1 is similar to modeling variations of C_0 w.r.t C_1 on any other $R_{i \neq 1}$. Thus, analysis of Local Conditional Probability could be restricted to single R_i at a time, as shown in Figure 8.7.

The different neighborhood systems \mathcal{N}^k , used in the RF models, $P^k(x|x_{\mathcal{N}^k})$, can be defined as (Refer Figure 8.7):

$$\mathcal{N}^k = \{C_{j|j \in \{|k|, 0^-, 0^+\}}\} \quad (8.8)$$

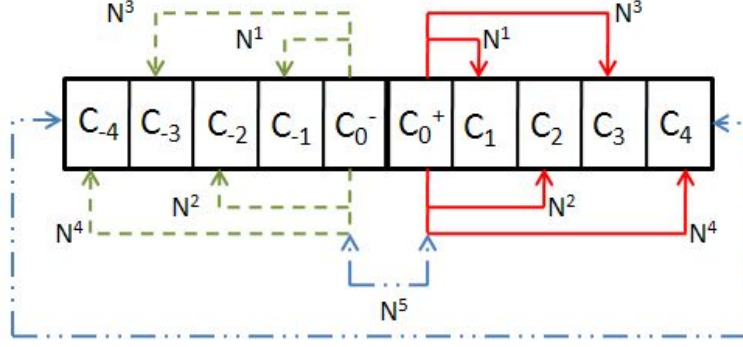


Figure 8.7: Neighborhood System.

8.3.2.4 Local Conditional Probability Density (LCPD)

To model the variations on the skin-region mask, we choose to build 2D histogram for each of the 5 RF over their unique neighborhood system. The design of the dimensions were such that they captured the various structural properties of true skin masks. The two dimensions (represented in a histogram pool \mathbf{H}^k) with individual element of the pool, \mathbf{z} , can be defined as:

- $\mathbf{H}^{k \in \{1,2,3,4\}} = \{\mathbf{z}\}$, where,

$$\mathbf{z} = [x_{C_{0^\pm}}, \delta(x_{C_{0^\pm}}, x_{C_{\pm k}})], \forall R_j \quad (8.9)$$

- $\mathbf{H}^{k=5} = \{\mathbf{z}\}$, where,

$$\mathbf{z} = [\mu(x_{C_{0^+}}, x_{C_{0^-}}), \mu(x_{C_{-4}}, x_{C_{+4}})], \forall R_j \quad (8.10)$$

where, x_{C_k} is the count of skin pixels in the block C_k . The two functions $\delta(\cdot, \cdot)$ and $\mu(\cdot, \cdot)$ are defined as

$$\delta(x_{C_{0^\pm}}, x_{C_{\pm i}}) = \begin{cases} x_{C_{0^+}} - x_{C_{+i}}, & i > 0 \\ x_{C_{-i}} - x_{C_{0^-}}, & i < 0 \end{cases} \quad (8.11)$$

$$\mu(a, b) = \frac{a+b}{2} \quad (8.12)$$

In order to estimate the LCPD on these 5 histogram pools, we use Parzen Window Density Estimation (PWDE) technique, similar to [242], with a 2D Gaussian window. Thus, each

of LCPD can now be defined as

$$P^k(\mathbf{z}) = \frac{1}{(2\pi)^{\frac{d}{2}} n h_{opt}^d} \sum_{j=1}^n \exp \left[-\frac{1}{2h_{opt}^2} \left(\mathbf{z} - \mathbf{H}_j^k \right)^T \Sigma^{-1} \left(\mathbf{z} - \mathbf{H}_j^k \right) \right]$$

where, n is the number of samples in the histogram pool \mathbf{H}^k , d is number of dimensions (in our case 2), Σ and h_{opt} are the covariance matrix over \mathbf{H}^k and the optimal window width, respectively, defined as:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}, \quad h_{opt} = \frac{\sigma_1 + \sigma_2}{2} \left\{ \frac{4}{n(2d+1)} \right\}^{1/(d+4)}$$

Figure 8.8 shows the 5 LCPDs learnt over a set of 390 training frontal face images.

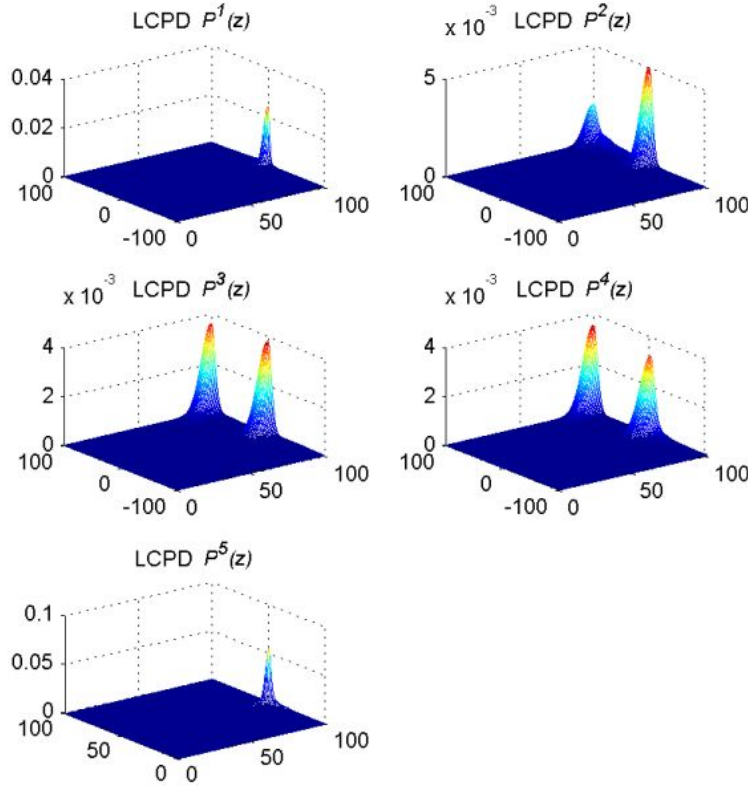


Figure 8.8: Frontal face Local Conditional Probability Density (LCPD) models.

8.3.2.5 Human Face Pose

During our studies we discovered that the structure of the skin-region varies based on the pose of detected face as shown in Figure 8.9. Combining face examples from different pose into one set of RFs seemed to dilute the LCPDs and hence the discriminating capability. This motivated us to design three different sets of RFs, one for each pose. This was accomplished by grouping true face detections into three piles, Turned right (r), Facing front (f), and, Turned Left (l).

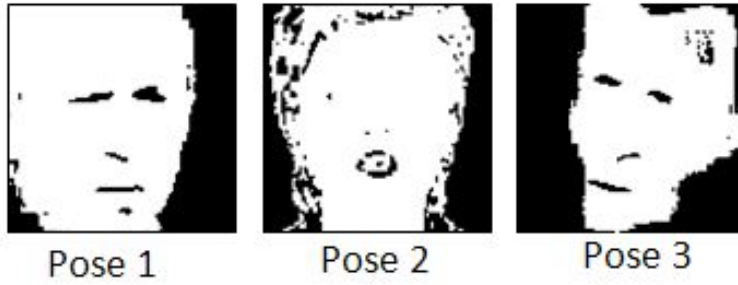


Figure 8.9: Skin-region masks.

Thus, the final set of LCPDs could be described by the super set.

$$P(\mathbf{z}) = \left\{ P_{m|m=\{r,f,l\}}^{k|k=\{1,\dots,5\}}(\mathbf{z}) \right\} \quad (8.13)$$

8.3.3 Combining Evidence

Given any test face detection output, \mathbf{z} is extracted (as described in Equation 8.9 and 8.10) and projected on the LCPD set $P(\mathbf{z})$ to get a set of likelihoods l_m^k . As in the case of any likelihood analysis, we combined the joint likelihood of multiple projections using log-likelihood function, $L_m^k = \ln(l_m^k)$, such that,

$$\prod_{\forall \mathbf{z} \in \mathbf{H}_m^k} \ln(l_m^k(\mathbf{z})) = \sum_{\forall \mathbf{z} \in \mathbf{H}_m^k} L_m^k(\mathbf{z}) \quad (8.14)$$

Given these log-likelihood values, one can set hard thresholds on each one of them to validate a face subimage discretely as *true* or *false*. We incorporated a piece-wise linear decision model (soft threshold) instead of a hard threshold on the acceptance of a face subimage.

This is illustrated in the Figure 8.10. Each LCPD $P^k(\mathbf{z})$ was provided with an upper and

lower threshold of acceptance and rejection respectively. The upper and lower bounds were obtained by observing $P^k(\mathbf{z})$ for the three face poses $P_{r,f,l}^k(\mathbf{z})$. Thus, any log-likelihood values lesser than the lower threshold (L_L) would result in a decision against the test input (Probability 0), while any log-likelihood value greater than the upper threshold (L_U) would be a certain accept (probability 1). Anything in between would be assigned a probability of acceptance. In order to combine the decisions from the five LCPD $P^k(\mathbf{Z})$, we resort to

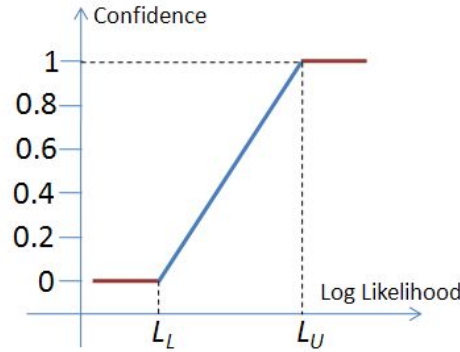


Figure 8.10: Soft threshold.

Dempster-Shafer Theory of Evidence.

8.3.3.1 Dempster-Shafer Theory of Evidence (DST)

The Dempster-Shafer theory is a mathematical theory of evidence [316] which is a generalization of probability theory with probabilities assigned to sets rather than single entities.

If X is an universal set with power set, $\mathbf{P}(X)$ (Power set is the set of all possible sub-sets of X , including the empty set \emptyset), then the theory of evidence assigns a belief mass to each subset of the power set through a function called the basic belief assignment (BBA), $m : \mathbf{P}(X) \rightarrow [0, 1]$, when it complies with the two axioms. a) $m(\emptyset) = 0$ and b) $\sum_{A \in \mathbf{P}(X)} m(A) = 1$. The mass, $m(A)$, of a given member of the power set expresses the proportion of all relevant and available evidence that supports the claim that the actual state belongs to A and to no particular subset of A . In our case, $m(A)$ correlates to the probability assigned by each of LCPDs towards the subimage being a face or not.

The true use of DST in our application becomes clear with the *rules* of combining

evidences which was proposed as an immediate extension of DST. According to the rule, the combined mass (evidence) of any two expert's opinions, m_1 and m_2 , can be represented as:

$$m_{1,2}(A) = \frac{1}{1-K} \sum_{B \cap C = A, A \neq \emptyset} m_1(B)m_2(C) \quad (8.15)$$

where,

$$K = \sum_{B \cup C = \emptyset} m_1(B)m_2(C) \quad (8.16)$$

is a measure of the conflict in the experts opinions. The normalization factor, $(1 - K)$, has the effect of completely ignoring conflict and attributing any mass associated with conflict to a null set.

The 5 LCPDs, $P^k(\mathbf{z})$, were considered as experts towards voting on the test input as a face or non-face. In order to use these mapped values in Equation 8.15 - 8.16, we normalized evidences generated by the experts to map between $[0, 1]$, and any conflict of opinions were added into the conflict factor, K . For the sake of clarity, we show an example of combining two expert opinions in Figure 8.11. The same idea could be extended to multiple experts.

		Expert 1's opinion	
		Face $m_1(B)$	Non-Face $m_1(C)$
Expert 2's Opinion	Face $m_2(B)$	Opinion Intersect $[m_1(B) * m_2(B)]$ (Sum in Numerator)	Opinion Conflict $[m_1(C) * m_2(B)]$ (Sum into K)
	Non-face $m_2(C)$	Opinion Conflict $[m_1(B) * m_2(C)]$ (Sum into K)	Opinion Intersect $[m_1(C) * m_2(C)]$ (Sum in Numerator)

Figure 8.11: An example of combining evidence from two experts under Dempster-Shafer Theory.

8.3.4 Coarse Pose estimation

Since the RF models were biased with pose information, we also investigated the possibility of determining the pose of the face based on the evidences obtained from the LCPDs. We noticed that the LCPDs $P^3(\mathbf{z})$, $P^4(\mathbf{z})$ and $P^5(\mathbf{z})$ were capable of not only discriminating faces from non-faces, but were also capable of voting towards one of 3 pose classes,

Looking right, Frontal, and Looking Left along with a confidence metric. Due to space constraints, the procedure is not explained in detail, but it is similar to what was followed for face versus non-face discrimination as explained in Section 8.3.3.

8.4 Testing the Abilities of the Face Detector

In all our experiments, Viola-Jones face detection algorithm [308] was used for extracting face subimages. The proposed face validation filter was tested on two face image data sets, 1. The FERET Color Face Database, and 2. An in-house face image database created from interview videos of famous personalities.

In order to prepare the data for processing, face detection was performed on all the images in both the data sets. The number of face detections do not directly correlate to the number of unique face images as there are plenty of false detections. We manually identified each and every face detection to be *true* or *false* so that ground truth could be established. The details of this manual labeling is shown below:

1. FERET

- Number of actual face images: 14,051
- Number of faces detected using Viola-Jones algorithm: 6,208
- Number of *true* detections: 4,420
- Number of *false* detections: 1,788 (28.8%)

2. In-house database

- Number of actual face images: 2,597
- Number of faces detected using Viola-Jones algorithm: 2,324
- Number of *true* detections: 2,074
- Number of *false* detections: 250 (10.7 %)

8.5 Results

In order to compare the performance of the proposed face validation filter, we defined four parameters:

1. Number of false detections (NFD)

$$\text{NFD} = \text{Count of false detections}$$

2. False detection rate (FDR):

$$\text{FDR} = \frac{\text{\# of false detections}}{\text{Total \# of face detections}} \times 100$$

3. Precision (P)

$$P = \frac{\text{\# of true detections}}{\text{\# of true detections} + \text{\# of false detections}}$$

4. Capacity (C)

$$C = \left(\frac{\text{\# of true detections}}{\text{\# of actual faces in database}} \right) - \text{FDR}$$

Table 8.1: Face detection validation results on FERET database.

	Before Validation	After Validation
NFD	1,788	208
FDR	28.8 %	3.35 %
P	0.7120	0.9551
C	0.026	0.281

Table 8.2: Face detection validation results on the in-house face database.

	Before Validation	After Validation
NFB	250	2
FDR	10.76 %	0.01 %
P	0.892	0.999
C	0.691	0.798

As explained in Section 8.3.4, the framework was extensible to perform coarse pose estimation. Figure 8.12 shows the result of passing two frames of a video sequence as input the face validation filter. The frames were extracted from a video of the same individual exhibiting arbitrary facial motion. The frames were 0.55 seconds apart. As can be noticed, the head pose is slightly different between the two frames. The pose estimation results are shown below the two frames.

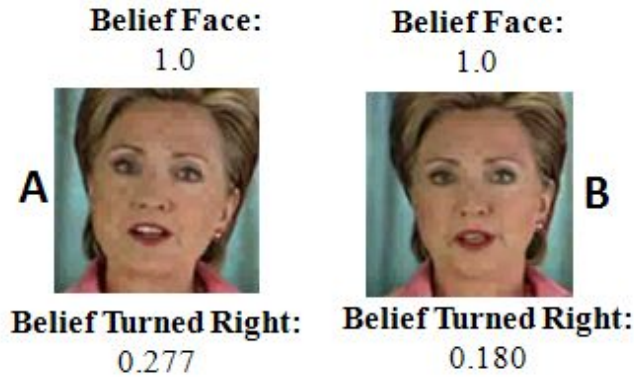


Figure 8.12: Coarse pose estimation.

8.6 Discussion of Results

Performance analysis of the proposed face validation filter can be understood through the four parameters defined in Section 8.5. **NFB** and **FDR** are direct measurements of the number of mistakes (naming non-faces as faces) made by the face detection algorithm on the two data sets. As can be verified from Table 8.1 and 8.2, there is a significant reduction in the false detections through the introduction of the filter.

The precision parameter, **P**, can be perceived as the probability that a face detection result retrieved at random will truly contain a face. It can be seen that the precision of the system drastically improves with the introduction of the face validation filter thereby assuring a *true* face subimage at the output.

The capacity parameter, **C**, measures the relative difference between face detection and false detection rates of a face detection system. Alternately, **C** can be considered to measure the net *true* face detection ability of any algorithm on a specific face data set. **C** ranges from -1 to 1 . -1 when none of the faces in the database are detected with all reported detections being wrong. 1 when all the faces in the database are detected with no false detections. It can be seen from Tables 8.1 and 8.2 that the capacity of the face detection system, when combined with face validation filter, is significantly higher and moves towards 1 . One can thus infer that the combined system has better *true* face detection ability.

Finally, Figure 8.12 shows the coarse pose estimation results. The two frames in the figure shows cases when the face is slightly turned right, with one (**A**) turned more right than the other (**B**). The face validation filter verifies that the faces are actually turned right and the belief values represent a scale on the amount of rotation. Since we did not do any specific mapping of the belief values to pose angle, we could not confirm quantitatively how accurate the pose estimations were. Through visual consort, one can verify that the labeling is meaningful.

SENSING DYNAMICS OF THE SOCIAL SCENE

As described earlier, determining the social scene and the dynamics of the social scene in front of the user requires that the person be detected in the camera video stream and they be tracked through in the video stream to determine the relative location of the individuals with respect to the user. That is, it is important to localize individuals in the video stream provided by the camera and track them through time. The problem of person localization in general is very broad in its scope and wide varieties of challenges such as variations in articulation, scale, clothing, partial appearances, occlusions, etc make this a complex problem. Narrowing the focus, this chapter targets person localization in real world video sequences captured from the camera of the social situational enrichment Assistant. Specifically, we focus on the task of localizing a person who is approaching the user to initiate a social interaction or just conversation. In this context, the problem of person localization can be constrained to the cases where the person of interest is facing the user.



Figure 9.1: Person of interest at a short distance from camera

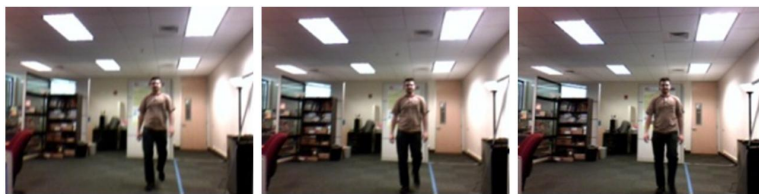


Figure 9.2: Person of interest at a large distance from camera

When such a person of interest is in close proximity, his/her presence can be detected by analyzing the incoming video stream for facial features (Figure 9.1), as explained in the previous chapter. But when such a person is approaching the user from a distance,

the size of the facial region in the video appears to be extremely small. In this case, relying on facial features alone would not suffice and there is a need to analyze the data for full body features (Figure 9.2). In this chapter, we focus on improving the effectiveness of the social interaction assistant by applying computer vision techniques to robustly localize people using full body features. Following section discusses some of the critical issues that are evident when performing person localization from the wearable camera setup of the SIA

9.1 Challenges in Person Localization from a wearable camera platform

A number of factors associated with the background, object, camera/object motion, etc. determine the complexity of the problem of person localization from a wearable camera platform. Following is a descriptive discussion of the imminent challenges that we encountered while processing the data using the SIA.

9.1.1 Background Properties

When the Social Interaction Assistant is used in natural settings, it is highly possible that there are objects in the background which move, thus causing the background to be dynamic. Also, there are bound to be regions in the background whose image features are highly similar to that of the person, thus leading to a cluttered background. Due to these factors, the problem of distinguishing the person of interest from the background becomes highly challenging in this context. Figures 9.3 and 9.4 illustrate the contrast in the data due to the nature of the background.

9.1.2 Object Properties

As we are interested in person localization, it can be clearly seen that the object is non-rigid in nature as there are appearance changes that occur throughout the sequence of images. Further, significant scale changes and deformities in the structure can also be observed. Also, when analyzing video frames of persons approaching the user, the basic image features in various sub-regions of the object vary vastly. For example, the image features from the facial region are considerably different from that of the torso region. Tracking detected

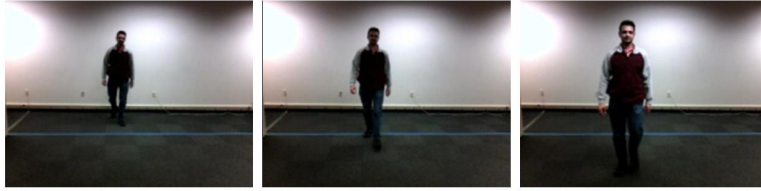


Figure 9.3: Simple Background



Figure 9.4: Complex Background

persons from one frame to another will require individualized tracking of each region to maintain confidence. This non-homogeneity of the object poses a major hurdle while applying localization algorithms and has not been studied much in the literature. Figure 9.5 shows the simplicity of the data when these problems are not present, while Figure 9.6 highlights complex data formulations in a typical interaction scenario.

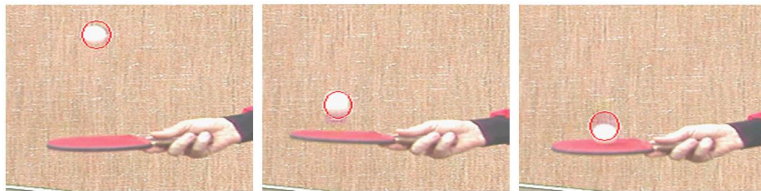


Figure 9.5: Rigid, Homogeneous Object



Figure 9.6: Non-Rigid, Deformable, Non-Homogeneous Object

9.1.3 Object/Camera Motion

Traditionally, most computer vision applications use a static camera where strong assumptions of motion continuity and temporal redundancy can be made. But in our problem, as

it is very natural for users to move their head continuously, the mobile nature of the platform causes abrupt motion in the image space (Compare Figure 9.7 and Figure 9.9). This is similar to the problem of working with low frame rate videos or the cases where the object exhibits abrupt movements. Recently, there has been an increase of interest in dealing with this issue in computer vision research [321] [322] [323] [324]. Some important applications which are required to meet real-time constraints, such as teleconferencing over low bandwidth networks, and cameras on low-power embedded systems, along with those which deal with abrupt object and camera motion like sports applications are becoming common place [324]. Though solutions have been suggested, person localization through low frame rate moving cameras still remains an active research topic.

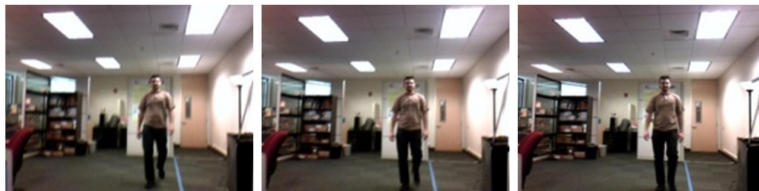


Figure 9.7: Static Camera

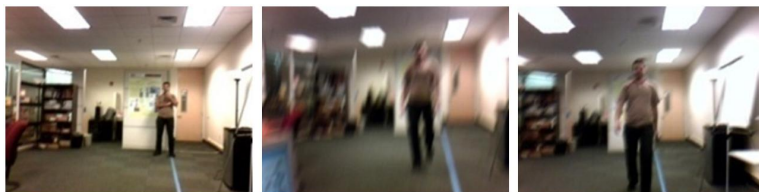


Figure 9.8: Mobile Camera

9.1.4 Other Important Factors Affecting Effective Person Tracking

As the SIA is intended to be used in uncontrolled environments, changing illumination conditions need to be taken into account. Further, partial occlusions, self occlusions, in-plane and out-of-plane rotations, pose changes, blur and various other factors can complicate the nature of the data. See Figure 9.9 for example situations where various factors can affect the video quality.

Given the nature of this problem, in this chapter we focus on the problem of robust localization of a single person approaching a user of the SIA using full-body features. Issues

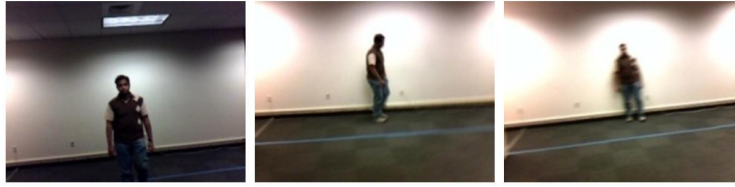


Figure 9.9: Changing Illumination, Pose Change and Blur

arising due to cluttered background along with object and camera motion have been handled towards providing robustness. In the following section we discuss some of the important related work in the computer vision literature.

9.2 Related Computer Vision Work in Person Localization and Tracking

Historically, two distinct approaches have been used for searching and localizing objects in videos. On one hand, there are detection algorithms which focus on locating an object in every frame using specific spatial features which are fine tuned for the object of interest. For example, haar-based rectangular features [325] and histograms of oriented gradients [326] can develop detectors that are very specific to objects in videos. On the other hand, there are tracking algorithms which trail an object using generic image features, once it is located, by exploiting the temporal redundancy in videos. Examples of features used by tracking algorithms include color histograms [327] and edge orientation histograms [328].

9.2.1 Detection Algorithms

As mentioned previously, detection algorithms exploit the specific, distinctive features of an object and apply learning algorithms to detect a general class of objects. They use information related to the relative feature positions, invariant structural features, characteristic patterns and appearances to locate objects within the gallery image. But, when the object is complex, like a person, it becomes difficult for these algorithms to achieve generality thereby failing even under minute non-rigidity. A number of human factors such as variations in articulation, pose, clothing, scale and partial occlusions make this problem very challenging.

When assumptions about the background cannot be made, learning algorithms which take advantage of the relative positions of body parts are used to build classifiers. The kind of low-level features generally used in this context are gradient strengths and gradient orientations [329] [326], entropy and haar-like features. Some of the well-known higher level descriptors are histogram of oriented gradients [326] and covariance features [330]. Efforts have been made to make these descriptors scale invariant as well.

In order to make these algorithms real-time, researchers have popularly resorted to two kinds of approaches. One category includes part-based approach such as Implicit Shape Models [321] and constellation models [331] which place emphasis on detecting parts of the object before integrating, while the other category of algorithms tries to search for relevant descriptors for the whole object in a cascaded manner[332]. Shape-based Chamfer matching [333] is a popular technique used in multiple ways for person detection as the silhouette gives a strong indication of the presence of a person. In recent times, Chamfer matching has been used extensively by the person detection and localization community. It has been applied with hierarchically arranged templates to obtain the initial candidate detection blocks so that they can be analyzed further by techniques such as segmentation, neural networks, etc. It has also been used as a validation tool to overcome ambiguities in detection results obtained by the Implicit Shape Model technique [334].

9.2.2 *Tracking Algorithms*

Assuming that there is temporal object redundancy in the incoming videos, many algorithms have been proposed to track objects over frames and build confidence as they go. Generally they make the simplifying assumption that the properties of the object depend only on its properties in the previous frame, i.e. the evolution of the object is a Markovian process of first order. Based on these assumptions, a number of deterministic as well as stochastic algorithms have been developed.

Deterministic algorithms usually apply iterative approaches to find the best estimate of the object in a particular image in the video sequence [332]. Optimal solutions

based on various similarity measures between the object template and regions in the current image, such as sum of squared differences (SSD), histogram-based distances, distances in eigenspace and other low dimensional projected spaces and conformity to particular object models, have been explored [332]. Mean Shift is a popular, efficient optimization-based tracking algorithm which has been widely used.

Stochastic algorithms use the state space approach of modeling dynamic systems and formulate tracking as a problem of probabilistic state estimation using noisy measurements [335]. In the context of visual object tracking, it is the problem of probabilistically estimating the object's properties such as its location, scale and orientation by efficiently looking for appropriate image features of the object. Most of these stochastic algorithms perform Bayesian filtering at each step for tracking, i.e. they predict the probable state distribution based on all the available information and then update their estimate according to the new observations. Kalman filtering is one such algorithm which fixes the type of the underlying system to be linear with Gaussian noise distributions and analytically gives an optimal estimate based on this assumption. As most tracking scenarios do not fit into this linear-Gaussian model and as analytic solutions for non-linear, non-Gaussian systems are not feasible, approximations to the underlying distribution are widely used from both parametric and non-parametric perspective.

Sequential monte-carlo based Particle Filtering techniques have gained a lot of attention recently. These techniques approximate the state distribution of the tracked object using a finite set of weighted samples using various features of the system. For visual object tracking, a number of features have been used to build different kinds of observation models, each of which have their own advantages and disadvantages. Color histograms [327], contours [336], appearance models, intensity gradients [337], region covariance, texture, edge-orientation histograms, haar-like rectangular features [332], to name a few. Apart from the kind of observation models used, this technique allows for variations in the filtering process itself. A lot of work has gone into adapting this algorithm to better perform in the context of visual object tracking.

While both the areas of detection and tracking have been explored extensively, there is an impending need to address some of the issues faced by low frame rate visual tracking of objects. Especially in the case of SIA, person localization in low frame rate video is of utmost importance. In this paper, we have attempted to modify the color histogram comparison based particle filtering algorithm to handle the complexities that occur mobile camera on the Social Interaction Assistant.

9.3 Conceptual Framework

As discussed in the previous section, detection and tracking offer distinctive advantages and disadvantages when it comes to localizing objects. In the case of SIA, thorough object detection is not possible in every frame due to the lack of computational power (on a wearable platform computing platform) and tracking is not always efficient due to the movement of the camera and the object's (interaction partner's) independent motion. Though there are clear advantages in applying these techniques individually, the strengths of both these approaches need to be combined in order to tackle the challenges posed by the complex setting of the SIA. In the past, a few researchers have approached the problem of tracking in low frame rate or abrupt videos by interjecting a standard particle filtering algorithm with independent object detectors [338]. In our experience, the Social Interaction Assistant offers a weak temporal redundancy in most cases. We exploit this information trickle between frames to get an approximate estimate of the object location by incorporating a deterministic object search while avoiding the explicit use of pre-trained detectors. Due to the flexibility in the design, particle filtering algorithms provide a good platform to address the issues arising due to complex data. These algorithms give an estimate of an object's position by discretely building the underlying distribution which determines the object's properties. But, real-time constraints impose limits on the number of particles and the strength of the observation models that can be used. This generally causes the final estimate to be noisy when conventional particle filtering approaches are applied. Unless the choice of the particles and the observation models fit the underlying data well, the estimate is likely to drift away as the tracking progresses. To mitigate these problems faced in the use of the SIA,

we propose a new particle filtering framework that gets an initial estimate of the person's location by spreading particles over a reasonably large area and then successively corrects the position through a deterministic search in a reduced search space. Termed as Structured Mode Searching Particle Filter (SMSPF), the algorithm uses color histogram comparison in the particle filtering framework at each step to get an initial estimate which is then corrected by applying a structured search based on gradient features and chamfer matching. The details of this algorithm are described in the next section.

9.4 Structured Mode Searching Particle Filter

Assuming that an independent person detection algorithm can initialize this tracking algorithm with the initial estimate of the person location, this particle filtering framework focuses on tracking a single person under the following circumstances, namely

- Image region with the person is non-rigid and non-homogeneous
- Image region with the person exhibits significant scale changes
- Image region with the person exhibits abrupt motions of small magnitude in the image space due to the movement of the camera.
- Background is cluttered.

The algorithm progresses by implementing two steps on each frame of the incoming video stream. In the first step (Figure 9.10), an approximate estimate of the person region is obtained by applying a color histogram based particle filtering step over a large search space. This is followed by a refining second step (Figure 9.11) where the estimate is corrected by applying a structured search based on gradient features and Chamfer matching. These two steps have been described in detail below.

9.4.1 Step 1: Particle Filtering Step

In the context of SIA, as the person of interest can exhibit abrupt motion changes in the image space, it is extremely difficult to model the placement of the person in the current



Figure 9.10: SMSPF - Step 1

image based on the previous frame's information alone. When such data is modeled in the Bayesian filtering based particle filtering framework, the state of each particle's position becomes independent of its state in the previous step. Thus, the prior distribution can be considered to be a uniform random distribution over the support region of the image.

$$p(x_t^i | x_{t-1}^i) = p(x_t^i) \quad (9.1)$$

As it is essential for particle filtering algorithm to choose a good set of particles, it would be useful to pick a good portion of them near the estimate in the previous step. By

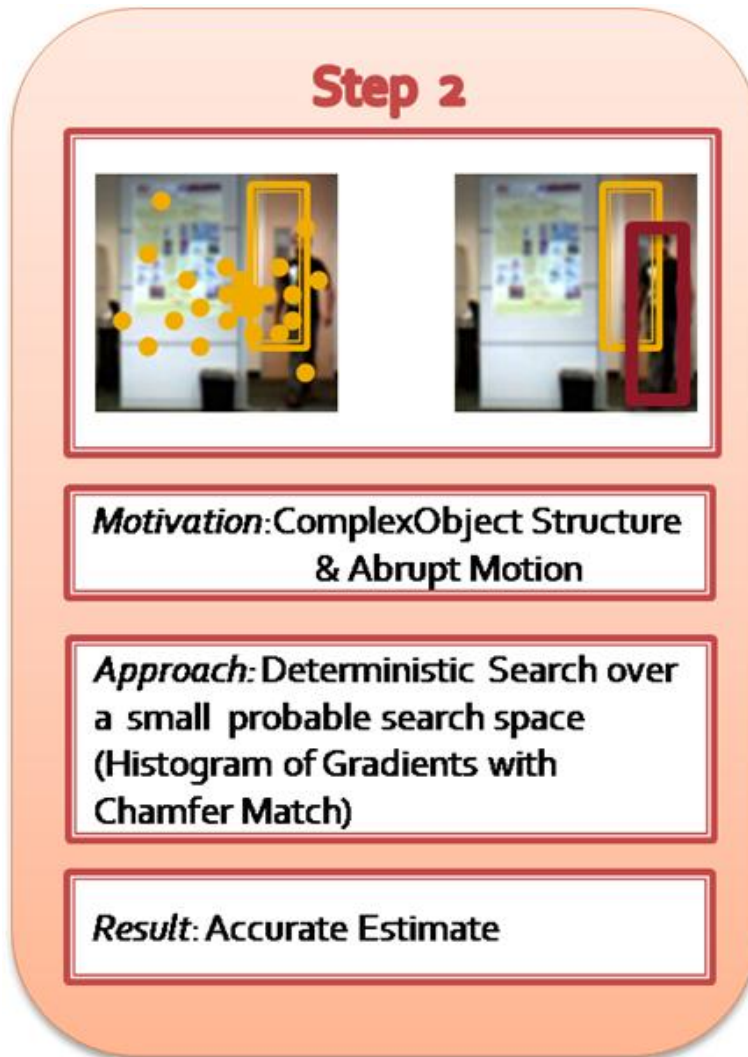


Figure 9.11: SMSPF - Step 2

approximating this previous estimate to be equivalent to a measurement of the image region with the person in the current step, the proposal distribution of each particle can be chosen to be dependent only on the current measurement

$$q(x_t^i | x_{t-1}^i Z_t) = q(x_t^i | Z_t) \quad (9.2)$$

Though the propagation of information through particles is lost by making such an assumption, it gives a better sampling of the underlying system. We employ a large variance Gaussian with its mean centered at the previous estimate for successive frame particle

propagation. By using such a set of particles, a larger area is covered, thus accounting for abrupt motion changes and a good portion of them are picked near the previous estimate, thus exploiting the weak temporal redundancy. As in [327], we have employed this technique using HSV color histogram comparison to get likelihoods at each of the particle locations. Since intensity is separated from chrominance in this color space, it is reasonably insensitive to illumination changes. We use an $8 \times 8 \times 4$ HSV binning thereby allowing lesser sensitivity to changes in V when compared to chrominance. The histograms are compared using the well-known Bhattacharyya Similarity Coefficient which guarantees near optimality and scale invariance.



Figure 9.12: Structured Search

With the above step alone, due to the small number of particles which are spread widely across the image, we can get an approximate location of the person. When such an estimate partially overlaps with the desired person region, the best match occurs between the intersection of the estimate and the actual person region as shown in Figure 9.12. But, it is not trivial to detect this partial presence due to the existence of background clutter. To handle this problem, we introduce a second step which uses efficient image feature representations of the desired person object and employs an efficient search around the estimate to accurately localize the person object.

9.4.2 Step 2: Structured Search

As the estimate obtained using widely spread particles gives the approximate location of the object, the search for the image block with a person in it can be restricted to a region around it. We have employed a grid-based approach to discretely search for the object of interest (a person) instead of checking at every pixel. By dividing the estimate into an $m \times n$ grid and sliding a window along the bins of the grid as shown in Figure 9.13, the search space can be restricted to a region close to the estimate. By finding the location which gives the best match with the person template, we can localize the person in the video sequence with better accuracy.

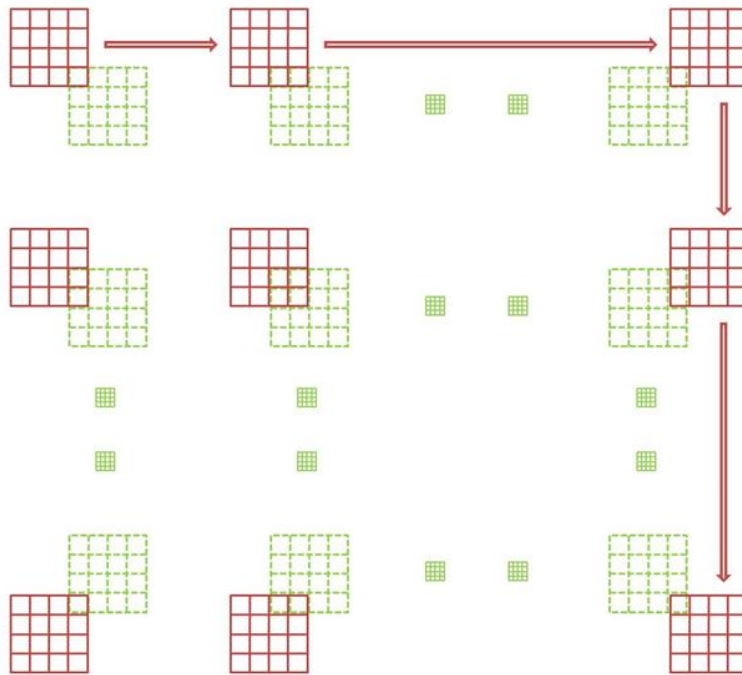


Figure 9.13: Sliding window of the Structured Search (Green: Estimate; Red: Sliding window).

If this search is performed based on scale-invariant features, then it can be extended to identify scale changes as well. In order to achieve search over scale, the estimate and the sliding window need to be divided into different number of bins. If the search is performed using smaller number of bins as compared to the estimate, then shrinking of the object can

be identified while searching with higher number of bins can account for dilation of the object. For example, if a $(m-1) \times (n-1)$ grid is used with the sliding window while a $m \times n$ grid is used with the estimate, then the best match will find a shrink in the object size. Similarly if an $m \times n$ grid sliding window is used with a $(m-1) \times (n-1)$ estimate grid, then dilations can be detected. It can be seen that this search is characterized by the number of bins $m \times n$ into which the sliding window and the estimate are divided. Based on the nature of the problem, the number of bins and the amount of sweep across scale and space can be adjusted. Currently, these parameters are being set manually, but the structured search framework can be extended to include online algorithms which can adapt the number of grid bins based on the evolution of the object.

If the object of interest was simple, then the best match across space and scale could be obtained by using simple feature matching techniques. But, due to the complex nature of the data, strong confidence is required while searching for the person region across scale. To this end, we propose to perform the structured search by analyzing the internal features of the person region as well as the external boundary/silhouette features and aggregating the confidence obtained from these two measures to refine the person location estimate in the image (Figure 9.14)

In literature, gradient based features have been widely used for person detection and tracking problems and their applicability has been strongly established by various algorithms like Histogram of Oriented Gradients (HoGs) [326]. Following this principle, we have used the Edge Orientation Histogram (EOH) features [328] in order to obtain the internal content information measure. For this purpose, a gradient histogram template (GHT) is initially built using a generic template image of a walking/standing person. This GHT is then compared with the gradient histogram of each structured search block using the Bhattacharyya histogram comparison as in [327] in order to find the block with the best internal confidence. In our implementation, orientations are computed using the Sobel operator and the gradients are then binned into 9 discrete bins. These features were extracted using the integral histogram concept [339] to facilitate computationally efficient searching.

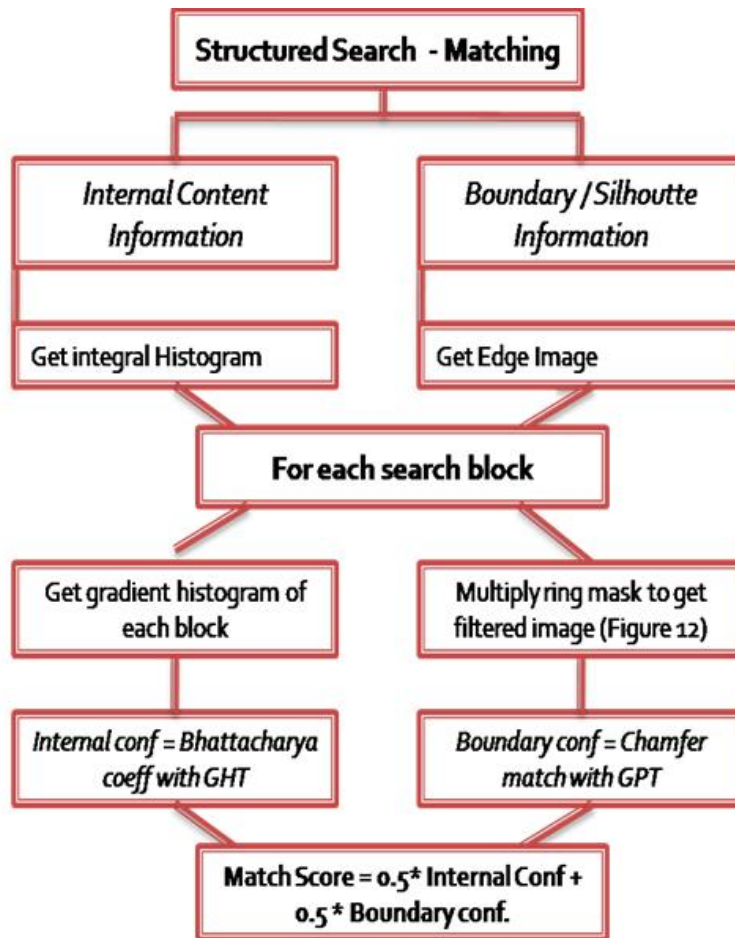


Figure 9.14: Structured Search Matching Technique

Similarly, in order to obtain the boundary confidence measure, a generic person silhouette template (GPT) (as shown in Figure 9.14) is used to perform a modified Chamfer match on each of the search blocks. In general, Chamfer matching is used to search for a particular contour model in an edge map by building a distance transformed image of the edge map. Each pixel value in a distance transformed image is proportional to the distance to its nearest edge pixel. In order to compare the edge map to the contour map, we convolve the edge image with the contour map. If the contour completely overlaps with the matching edge region, we get a chamfer match value of zero. Based on how different the edge map is to the template contour, the chamfer match score will increase and move towards 1. A chamfer match score of 1 implies a very bad match.

While the theory of chamfer matching offers elegant search score, in reality, especially with clutter within the object's silhouette, it is very difficult to get an exact match score. In SIA, since the data is very noisy and complex, certain modifications need to be made with the Chamfer matching algorithm in order to achieve good performance. The following section details a modified Chamfer match algorithm introduced in this work.

9.4.3 *Chamfer Matching in Structured Search*

As discussed above, Chamfer matching gives a measure of confidence on the presence of the person within an image based on silhouette information. We have incorporated this confidence into the structured search in order to detect the precise location of the person around the particle filter estimate. An edge map of the image under consideration is first obtained which is then divided into $(m \times n)$ windows in accordance with the structured search and an elliptical ring mask is then applied to each of these windows as shown in Figure 9.15. This mask is applied so as to eliminate the edges that arise due to clothing and background thereby emphasizing the silhouette edges which are likely to appear in the ring region if a window is precisely placed on the object perimeter. A distance transformed image of the window is then obtained using the masked edges.

By applying the modified chamfer matching (with a generic person contour resized to the current particle filter estimate), a confidence number in locating the desired object within the image region can be obtained. Similar to the Chamfer matching as before, a value close to 0 indicates a strong confidence of the presence of a person and vice versa. As 1 is the maximum value that can be obtained by the chamfer match, this measure can be incorporated into the match score of the structured search using the following equation.

$$\text{BoundaryConf} = (1 - \text{ChamferMatch}) \quad (9.3)$$

The standard form of Chamfer Matching gives a continuous measure of confidence in locating an object in an edge map. But, in our case, when the elliptical ring mask is used to filter out the noisy edges in each search block, this nature of Chamfer match is lost. Since

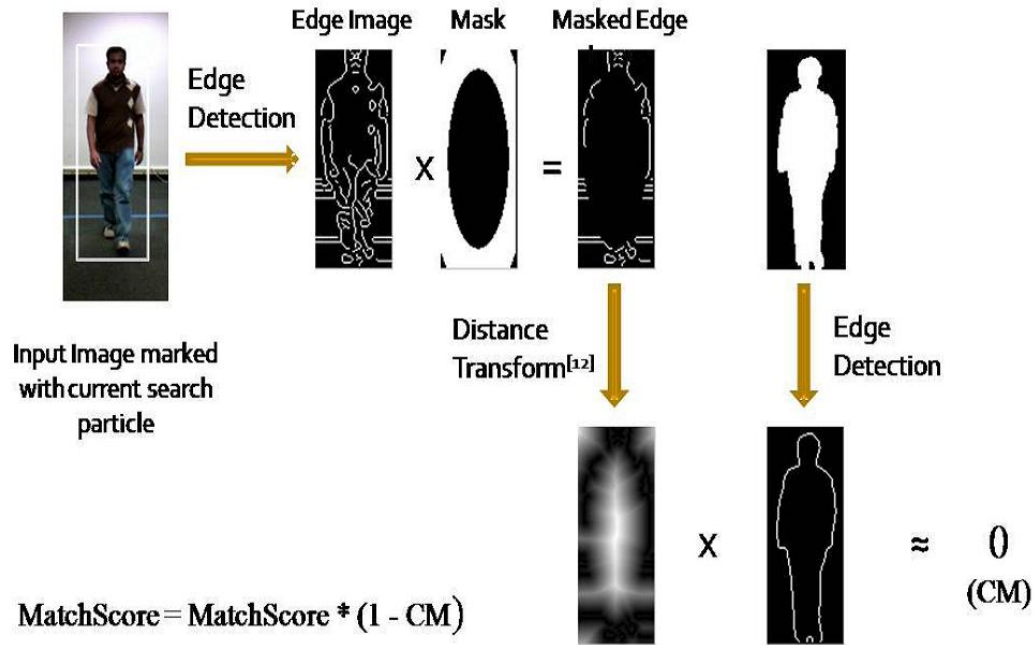


Figure 9.15: Incorporating Chamfer Matching into Structured Search

the primary goal of the structured search is to find a single best matching location of the person, it is more advantageous to use the filter mask at the cost of losing this continuous nature of the chamfer match. Further, as it is very likely that the person region is close to the approximate estimate obtained from the first step, one of the search windows of the structured search is bound to capture the entire person object thus resulting in a good match score.

From the above discussion, it can be seen that combining the knowledge about the internal structure of the person region with the silhouette information results in a greater confidence in the SMSPF algorithm. Further, using such complementary features in the structured search robustly corrects the approximate estimate obtained from the particle filtering step while handling various problems associated with search across scale.

9.5 Experiments and Datasets

9.5.1 Datasets

The performance of the structured mode searching particle filter (SMSPF) has been tested using three datasets where a single person faces the camera while approaching it. There are significant scale changes in each of these sequences. Further, non-rigidity and deformability of the person region can also be clearly observed. Different scenarios with varying degrees of complexity of the background and camera movement have been considered. Following is a brief description of these datasets.

- (a) *DataSet*¹: Plain Background; Static Camera; 320x240 resolution
- (b) *DataSet*²: Slightly cluttered Background; Static Camera; 320x240 resolution
- (c) *DataSet*⁴: Cluttered Background; Mobile Camera; 320x240 resolution

Figure 9.16 shows the sample results on each of the datasets used.

9.5.2 Evaluation Metrics

In order to test the robustness of this algorithm and the applicability in complex situations, its performance has been compared with the Color Particle Filtering algorithm [333]. Assuming that a detection algorithm can detect persons in at least some frames, the image region containing the person in each of the test sequences has been manually set. The following two criteria have been used to evaluate their performance,

- Area Overlap (A0)
- Distance between Centroids (DC)

¹Collected at CUBiC

²CASIA Gait Dataset B with subject approaching the camera ³

⁴Collected at CUBiC



(a) SMSPF Results on a sequence from Dataset 1



(b) SMSPF Results on a sequence from Dataset 2



(c) SMSPF Results on a sequence from Dataset 3

Figure 9.16: SMSPF Results

Manually labeled rectangular regions around the person in the image have been used as the ground truth. Suppose $gTruth_i$ is the ground truth in the i^{th} frame and $track_i$ is the rectangular region output by a tracking algorithm, then the area overlap criterion is defined as follows

$$AO(gTruth_i, track_i) = \frac{Area(gTruth_i \cap track_i)}{AO(gTruth_i \cup track_i)} \quad (9.4)$$

The average area overlap can be computed for each data sequence as

$$AvgAOR = \frac{1}{N} \sum_{i=1}^N AO \quad (9.5)$$

Similar to [340], we use Object Tracking Error (OTE) which is the average distance between the centroid of the ground truth bounding box and the centroid of the result given by a tracking algorithm

$$OTE = \frac{1}{N} \sum_{i=1}^N \sqrt{(Centroid_{gTruth_i} - Centroid_{Truth_i})} \quad (9.6)$$

In order to evaluate the performance of these algorithms using a single metric which encodes information from both area overlap and the distance between centroids, we have used a measure termed as the Tracking Evaluation Measure (TEM) which is the harmonic mean of the average area overlap fraction (AvgAOR) and a non-linear mapping of the Object tracking error (OTE).

$$TEM = 2 * \frac{AvgAOR.e^{-k.OTE}}{AvgAOR + e^{-k.OTE}} \quad (9.7)$$

where k is a constant which exponentially penalizes the cases where the distance between centroids is large.

9.6 Results

Particle Filtering has been widely used to handle complex scenarios by maintaining multiple hypotheses. As mentioned in [336], in order to handle abrupt motion changes, it is essential that the particles are widely spread while tracking. Following this principle, we have compared the performance of color particle filter (PF) [333] and the structured mode searching particle filter (SMSPF) by using a 2-D Gaussian with large variance as the system model. The position of the person and its scale have been included in the state vector. In order to compensate for the computational cost of structured search, only 50 particles were used for the SMSPF algorithm while 100 particles were used for the PF algorithm. A 10x10 grid with a sweep of 8 steps along the spatial dimension and 3 steps along the scale dimension were incorporated in the structured search.

Figure 9.17 and Figure 9.18 illustrate the comparison of the area overlap ratio and the distance between centroids at each frame of an example sequence. The sample frames are shown beside the tracking results. From Figure 9.17(a), it is evident that the SMSPF algorithm (red) shows a significant improvement over the color particle filter algorithm

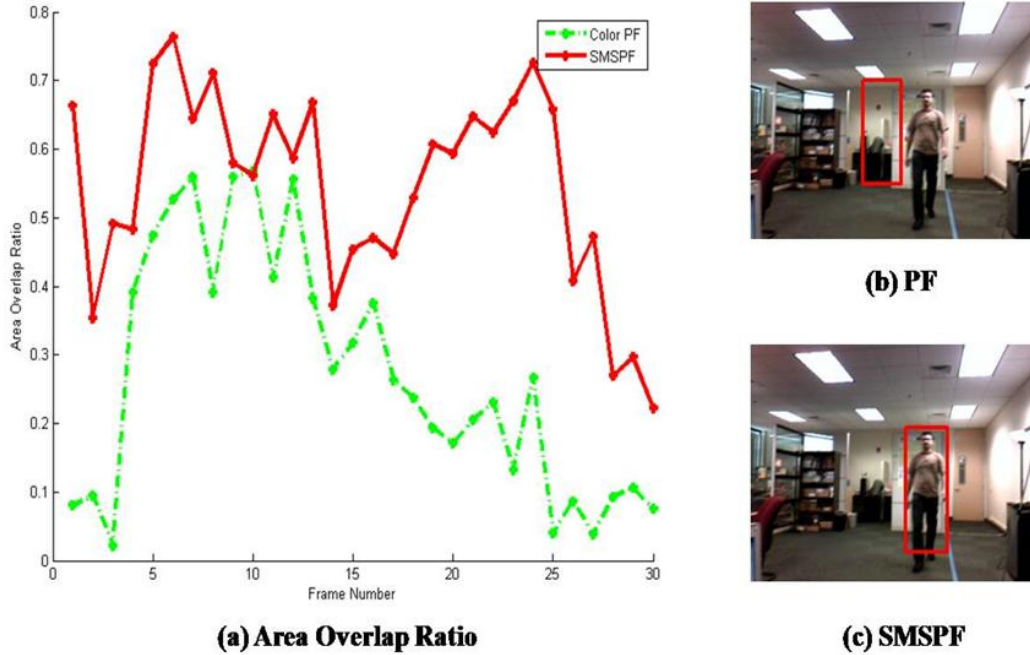


Figure 9.17: AO (Dotted Line: Color PF; Solid Line: SMSPF)

(green). Here, the area overlap ratio using SMSPF is much closer to 1 in most of the frames while the color particle filter drifts away causing this measure to be closer to 0. The distance between centroids measure also indicates a greater precision of the SMSPF algorithm as seen in Figure 9.18(a) where the distance between centroids using color particle filter is much higher than that with SMSPF (≈ 0).

Figure 9.19, Figure 9.20 and Figure 9.21 show the Tracking Evaluation Measure (TEM) for Datasets 1, 2 and 3. In majority of the cases, the SMSPF algorithm outperforms the color particle filtering algorithm with a higher TEM score.

The results presented as a comparison between Color PF and SMSPF shows that incorporating a deterministic structured search into the stochastic particle filtering framework improves the person tracking performance in complex scenarios. The SMSPF algorithm strikes a balance between specificity and generality offered by detection and tracking algorithms as discussed in Section 2. It uses specific structure-aware features in the search in order to handle non-homogeneity of the object and the cluttered nature of the background. On the other hand, generality is maintained by using simple, global features in the parti-

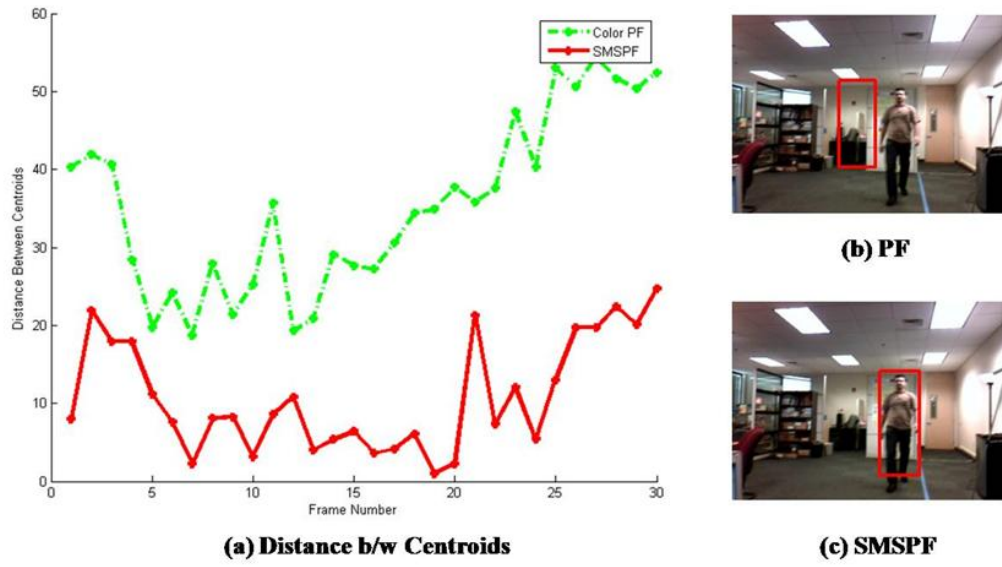


Figure 9.18: DC(Dotted Line: Color PF; Solid Line: SMSPF)

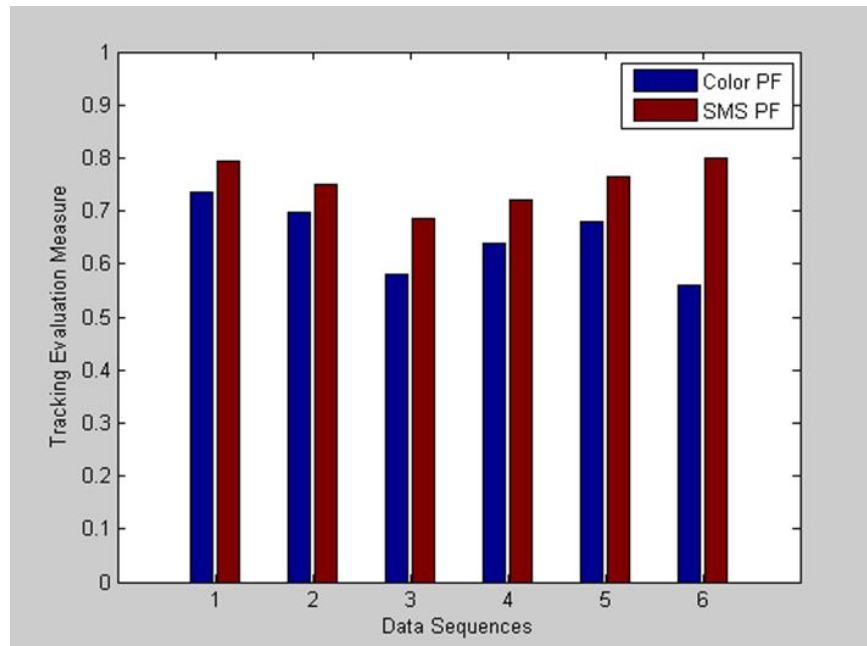


Figure 9.19: Evaluation Measure for DataSet 1

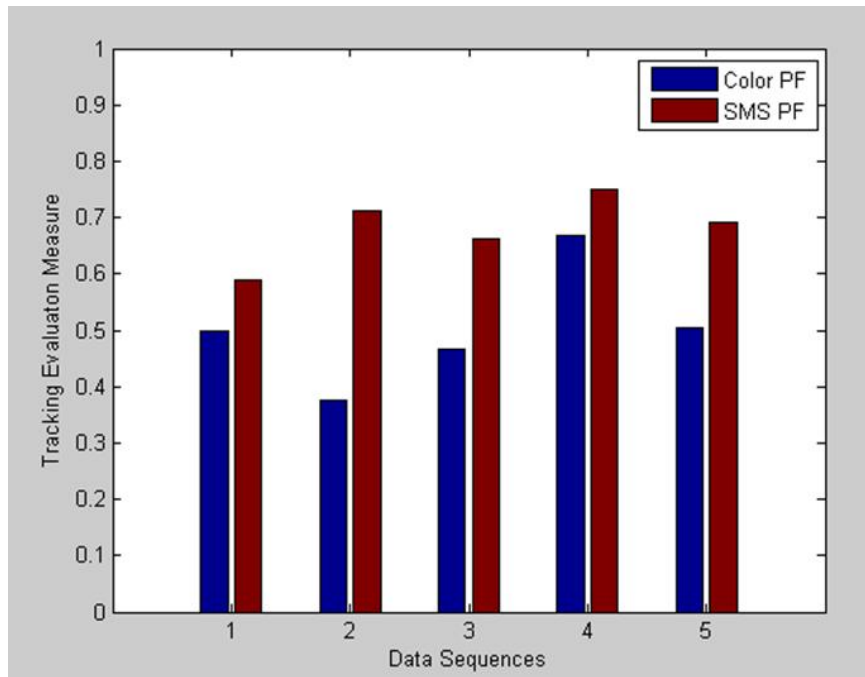


Figure 9.20: Evaluation Measure for DataSet 2

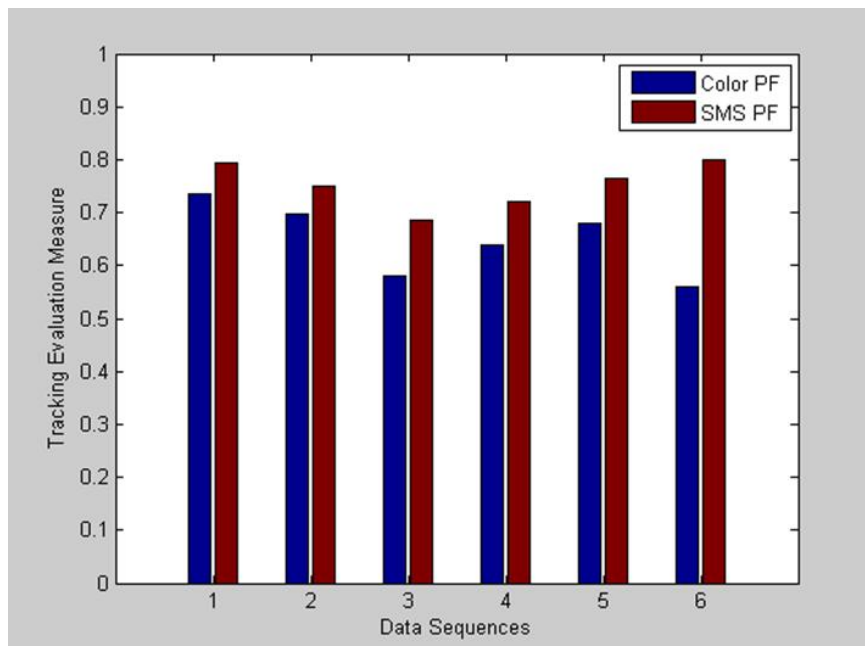


Figure 9.21: Evaluation Measure for DataSet 3

cle filtering framework so as to handle non-rigidity and deformability of the object. The clear advantage of using the structured search can be observed on the complex Dataset 3 which encompasses most of the challenges generally encountered while using the Social Interaction Assistant.

COMMUNICATING SOCIAL SCENE DYNAMICS

Like facial expression data extracted from a camera, the social scene data (group interaction data) has a very high bandwidth and communicating this data to individuals who are blind, requires that the data be modulated on a modality that does not interfere with the user's cognitive capacity. As described earlier, the use of audio cueing in such contexts may not be of high value as it has the tendency to overload the user's natural capacity to sense their environment. To this end, as before, we resort to the use of haptic technologies as augmented interface. Tables 6.1 through 6.4 introduced a plethora of haptic technologies used in communicating remote interpersonal data. While these technologies provided means of communicating specific interpersonal data, no one technology provided means of communicating distances and direction information to a user. To this end, we explored haptic technologies specifically designed for communicating spatial orientations. This chapter describes an alternative delivery modality: a vibrotactile belt that can convey non-verbal communication cues to individuals who are blind or visually impaired. Specifically, we focus on the non-verbal cue listed in Section 3.3.1 of Chapter 3: helping users perceive the number of people in their visual field, and the relative direction and distance of each individual with respect to the user. In some social situations, location information is available through audible cues, but this is not always the case. For example, when a group of friends approaches all of them may smile but only some may offer a verbal greeting, or a passing co-worker may nod to you in the hallway without exchanging a verbal greeting. These non-verbal communications are common occurrences, but are not accessible to someone who is blind.

10.1 Proposed Framework of Social Scene Structure Delivery

As shown in Figure 10.1, the output of the face detection/tracking process (indicated by a green rectangle on the image) provided by the camera on the Social Interaction Assistant is directly coupled with a vibrotactile haptic belt. Every frame in the video sequence captured

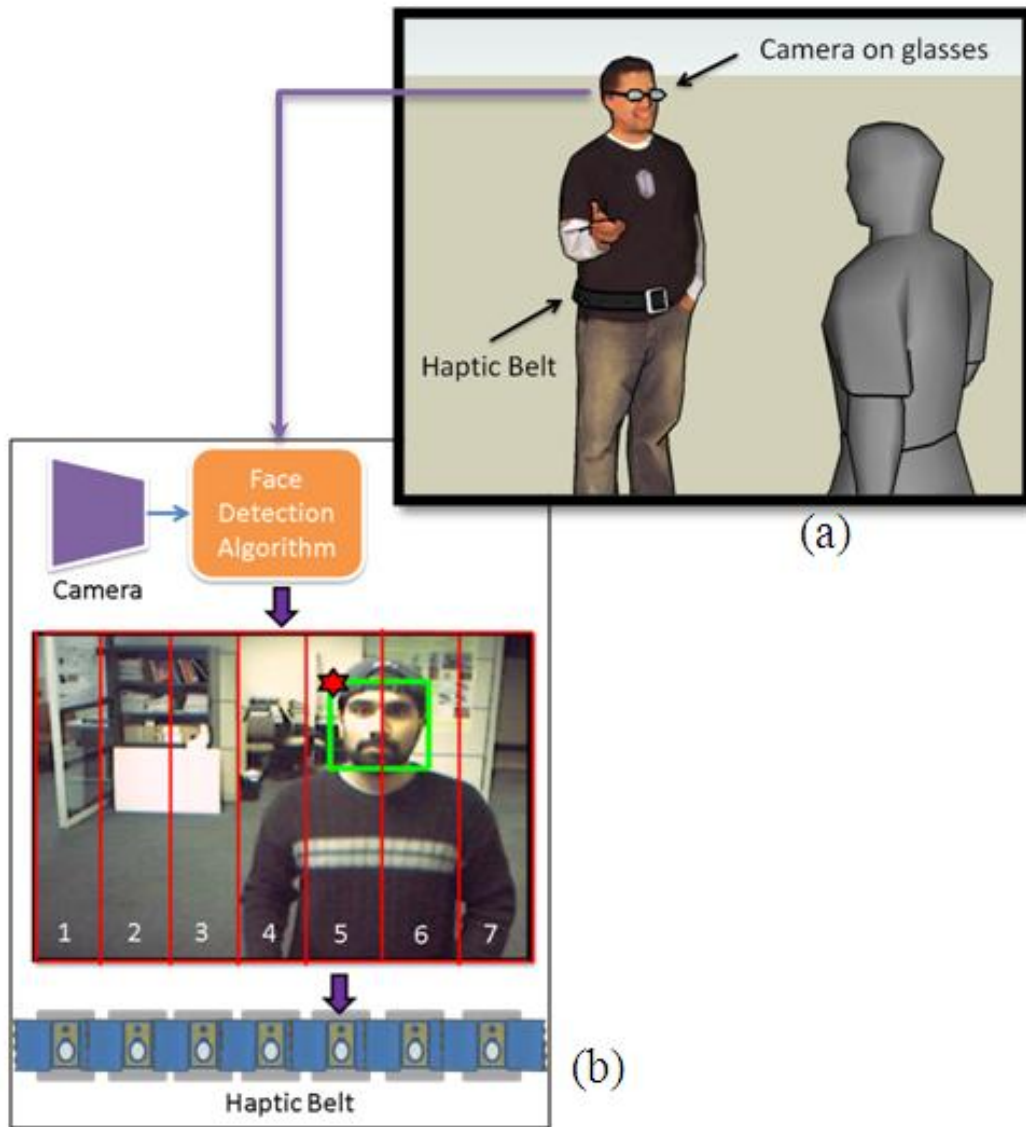


Figure 10.1: (a) Typical use of the social interaction assistant, a third person perspective on the use case scenario, (b) An example of face detection being translated to vibrations on the haptic belt.

by the Social Interaction Assistant is divided into N regions equal to N vibrotactile actuators on the haptic belt. In the example shown in the figure, there were 7 vibrators on the belt. After face detection, the region to which the top-left corner of the face detection output belongs is identified (as shown by the star in Figure 3). This region directly corresponds to the tactor on the belt that needs to be activated to indicate the direction of the person with respect to the user. The duration of a vibration indicates the distance between the user and the person in his or her visual field. The longer the vibration, the closer the people are, which is estimated by the face image size determined during the face detection process.

An overall perspective of the system and its process flow is given below. When a user encounters a person in his or her field of view, the face is detected and recognized (if the person is not in the face database, the user can add it). The delivery of information comprises two steps: Firstly, the identity of the person is audibly communicated to the user (we are currently investigating the use of tactons [341] to convey identities through touch, but this is part of future work). Secondly, the location of the person is conveyed through a vibrotactile cue in the haptic belt, where the location of the vibration indicates the direction of the person and the duration of vibration indicates the distance between the person and the user. Based on user preference, this information can be repeatedly conveyed with every captured frame, or just when the direction or distance of the person has changed. The presence of multiple people in the visual field is not problematic as long as faces are not occluded and can be detected and recognized by the Social Interaction Assistant. We are currently investigating how to effectively and efficiently communicate non-verbal communication cues when the user is interacting with more than one person.

10.2 Related Work in Haptic Vibrotactile Technology for Information Delivery

Of all the modalities that engage the human somatosensory system, vibrotactile stimulation has become very popular in the recent past due to the sophistication and unobtrusiveness of vibrotactile displays [342], as well as their portability and wearability [343]. But this new modality is far from displacing the primary delivery modalities due to the fact that haptics (touch) is a low bandwidth channel compared to audio or video. Previously, complex vi-

bratory pulses have been designed using combinations of vibration dimensions [344][345], such as vibration frequency, amplitude, duration, rhythm and location, and by using human psychophysical perception (like sensory saltation [346]). There are infinite ways to map meanings to vibration dimensions, but conceptually, there are two extremes: symbolic and literal. On one end of the spectrum, tactons [341], or tactile icons, use a symbolic mapping to arbitrarily assign meaning to vibration dimensions. On the other end, a literal mapping assigns vibrotactile cues to intuitive somatosensory signals that humans are already acquainted with, such as a shoulder tap to obtain one's attention. Encoding schemes may also fall somewhere in between such that the vibrotactile cues may be intuitive, but still require training. Studies on symbolic and literal mappings have shown an extraordinary increase in information delivery bandwidth for vibratory cues, thereby making a case for vibrotactile stimulation as a potential alternative (or at least an augmentation) to audio and video.

Vibrotactile displays have been implemented in a variety of form factors including desktop displays, handheld devices, and wearable systems, such as gloves [257], jackets [347], and jewelry [263]. In this paper we focus our discussion to vibrotactile displays worn around the waist, commonly referred to in the literature as haptic or vibrotactile belts. Vibrotactile belts have found a number of applications including, but not limited to, pedestrian navigation [348][349] [350] (vibrations guide users from a starting point to their destination); balance control [351] for people with vestibular damage (vibrations convey tilt information); virtual reality [352] (vibrations indicate collisions with virtual objects); spatial orientation aids for pilots [353] and astronauts [354] (vibrations provide spatial orientation towards magnetic north or Earth's gravity vector in zero-gravity environments); psychophysical study of human vibrotactile perception [345] [355] (experiments on vibrotactile spatial acuity, spatio-temporal pattern perception, saltation, etc.); and social interaction assistant aids [252] for individuals who are blind or visually impaired (vibrations are used to communicate nonverbal cues). Unlike other form factors, belts tend to be physically discreet and part of almost all everyday clothing. A variety of vibrotactile belt designs and implementations have been proposed in the literature (please refer to Section II for a detailed analysis). However, existing designs have two primary limitations: (1) Limited applicabil-

ity due to application-specific designs; and (2) Usability and performance requirements tend to be secondary to functionality, thereby forcing readers to question the real-world use of the application itself. This is the natural inclination of a technology-centric, as opposed to human-centric, approach towards interface design. A human-centered design strategy critically accounts for all users of the technology throughout the lifecycle of the design and development of a human machine interface. In this work, we generalize the scope of our users to include both customers: end users of a specific technology; and developers (engineers, scientists, and researchers): those modifying the product for novel applications.

Our literature survey revealed over twenty vibrotactile belt designs from academic publications and electronics hobby forums. We've selected a subset for discussion here based on the maturity of their implementation and availability of information regarding implementation details.

Cholewiak et al. [355] introduced a reconfigurable and scalable haptic belt design for use in human haptic perception experiments, where vibration motors were wired directly to a waveform generator, and attached via Velcro onto an elastic belt. The belt was specifically intended for psychophysical experiments, and its wired implementation limits portability, ease of movement, unobtrusiveness and discreetness. Van Erp et al. [348] presented a wireless, elastic vibrotactile belt for waypoint navigation. The belt consisted of eight vibration motors with adjustable locations. The belt was controlled by a minicomputer placed inside a backpack worn by the user. The paper provides no information regarding the scalability of the belt, i.e., the option of removing or adding vibration motors. Moreover, it is unclear if the amplitude and/or frequency of the vibrations can be adjusted. For studying human haptic perception, Jones and Ray [345] built a wireless haptic belt made of fabric consisting of 8 vibration motors held by Velcro. A back display was also constructed, which consisted of a four-by-four matrix of vibration motors. The locations of the vibration motors were adjustable, but the paper does not mention whether amplitude or timing could be controlled, nor is there any mention of the capability to add or remove vibration motors. Further, the bulkiness of the system, and its excessive cabling, could limit ease of

movement, unobtrusiveness and discreetness.

ActiveBelt [350] is a wireless haptic belt for pedestrian navigation, among other applications. The belt consisted of eight fixed vibration modules with elastic between vibration sites and used a large, onboard processing unit. Dimensions of the vibratory signals, such as frequency and timing, could be altered, but the reconfigurability and scalability of ActiveBelt is limited given its fixed vibration motors. Further, although the paper claims universal accessibility in that the belt can adapt to varying waist sizes, this may be only partially true—from our own past experiences, extreme waist sizes (either very small or very large) may not be able to use such an implementation.

Ferscha et al. [356] presented a wireless vibrotactile belt for spatial awareness. Vibratory dimensions, such as intensity and timing, could be altered in a portable and lightweight design. However, since the belt used eight fixed vibration motors, its reconfigurability and scalability is limited. The Tactile Wayfinder [349], by Heuten et al., is a wireless vibrotactile belt for pedestrian navigation. It has many of the same advantages and disadvantages of Ferscha et al.'s belt design, but with a few differences; one advantage being The Tactile Wayfinder has an available API for application creation.

Perhaps the most accomplished of the aforementioned belt designs is the TactaBelt by Lindeman et al. [352], which consisted of eight vibration motors connected via Velcro to neoprene. The vibration motors of the belt are reconfigurable and scalable, and their vibratory dimensions are adjustable. Although the TactaBelt is functional and rich in features, there is little to no discussion regarding the usability and performance of the belt—this is also a reoccurring problem with all aforementioned belt designs. The rigidity and durability of this belt is questionable given that vibration motors were attached to the belt via Velcro. Whereas this solution may work in controlled environments, such as a virtual reality setup in a laboratory, it's unlikely to work well in real world conditions and under everyday use.

In [357], Ram and Sharf introduced The People Sensor: an electronic travel aid, for individuals who are blind, designed to help detect and localize people and objects in front of the user. The distance between the user and an obstacle is found using ultrasonic sensors

Table 10.1: Design Requirements for Vibrotactile Belts

Usability	Functionality	Performance
<i>Limited Cumber</i>	<i>Expressiveness</i>	✓ Robustness and rigidity
✓ Easy to take on/off	✓ Dimensions of vibrations changeable	✓ Reliable
✓ Doesn't hinder movement	<i>Scalability</i>	✓ Long wireless communication range
✓ Comfortable	✓ Tactors can be added or removed	✓ Negligible latency in wireless communication
✓ Ergonomic	<i>Reconfigurability</i>	✓ Long battery life
✓ Unobtrusive	✓ Position of tactors can be changed	✓ Rechargeable or replaceable batteries
✓ Lightweight	✓ API is available	
✓ Adaptive	<i>Portability</i>	
<i>Intuitiveness</i>	✓ Wearable	
✓ Easy to learn and use	✓ Wireless	
<i>Discreetness</i>		
✓ Physically discreet		
✓ Silent		

and communicated through the rate of short vibratory pulses, where the rate is inversely proportional to distance. However, the researchers did not do any user testing to determine the usefulness of their technology.

10.3 Design Requirements

Identifying the shortcomings of vibrotactile belt designs, reviewing existing design guidelines in the literature, and combining these with our own past experiences, we've compiled a set of design requirements for vibrotactile belts, depicted in Table I.

In the above table, usability is the most important metric that captures the capability of a haptic platform to be used for exploring novel applications; in other words, if there are usability issues in a research platform, it will bias the outcome of any research experiment, thereby distracting the researcher from the true outcomes of an experiment. Following usability, functionality takes the next higher precedence, as it allows a researcher to configure the device to his or her novel application needs. Offering higher functionality allows adaptability of the research platform to various experiments. Finally, performance

captures the lenience offered by the platform during experimental use. Mostly, higher performance reduces the researcher's requirement to focus attention on the research platform, and allows him/her to focus on the study itself. We discuss some of the existing work in eliciting such design requirements for vibrotactile belt and add design considerations that we have identified through our experience.

Regarding the usability of a vibrotactile belt, Lindeman et al. [352] described a vibrotactile wearable device with limited cumber as one that is easy to put on or take off, and doesn't hinder movement with excessive wiring and bulky modules. Adding to this description of limited cumber, we include factors such as comfort and unobtrusiveness [358], ergonomics, lightweight and adaptability to fit different waist sizes. A vibrotactile belt should be intuitive so that it is easy to learn to use from both an end-user's perspective, as well as a developer's perspective. The latter will have much more vested interest in reconfiguring the belt for his or her intended application. Lastly, a vibrotactile belt should be discreet in that it is physically discreet and silent. As belts are a common part of everyday attire, keeping the design of vibrotactile belts close to accustomed dressing attire will help gain wider acceptance among users. Vibration motors can be noisy, which when used in public, can be distracting to those around us. Hence, vibrotactile modules should be designed to reduce noise.

Lindeman et al. [352] proposed three functionality attributes: expressiveness, scalability and reconfigurability as being important for a vibrotactile display. The first attribute, expressiveness, was met by providing variability of vibration dimensions: intensity, timing and location. However, the paper gives little detail about what exactly defines scalability and reconfigurability of a vibrotactile belt. We extend their work to define scalability as the capability to add/remove factors to/from a vibrotactile belt without performance degradation; and reconfigurability, which is related to the adaptability of the belt to different applications and uses, is defined as the capability to (1) easily change the placement of factors on a vibrotactile belt, and (2) easily change the vibrotactile belt's functions through an Application Programming Interface (API). Lastly, portability is an important functionality

influenced by its wearability and wireless connectivity. Attributes that describe performance design requirements include durability, long wireless communication range, negligible latencies in wireless communication, long battery life and replaceable/rechargeable batteries. Although the importance of these attributes will largely depend on an application's minimum performance requirements, it's recommended that all of the proposed attributes be taken into account when developing a versatile vibrotactile belt.

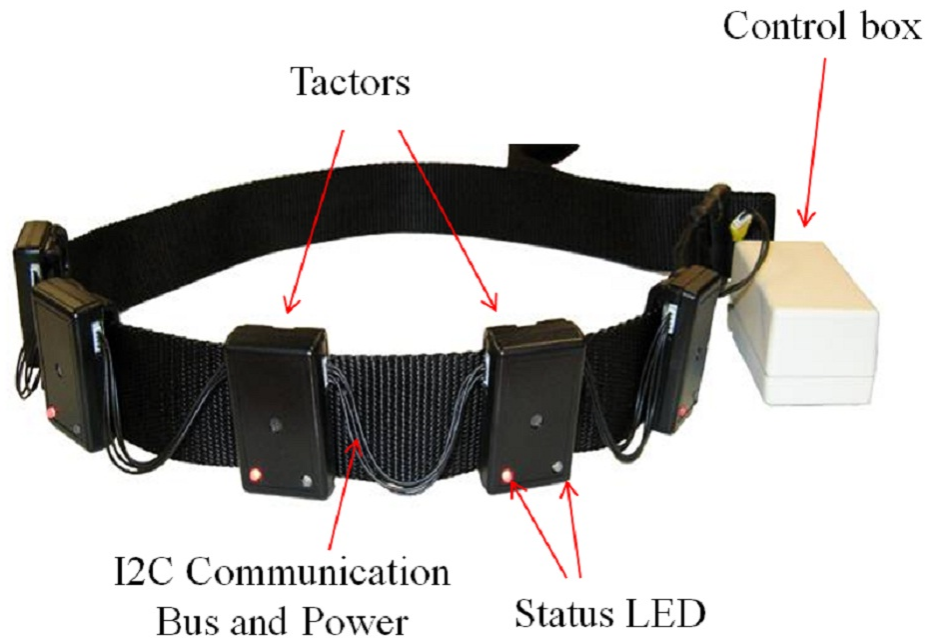
As mentioned earlier, the versatility and usefulness of existing vibrotactile belt implementations are severely limited due to an application-specific focus. Such a non-structured approach results in replication of work between researchers and developers. Our goal, through this paper, is to establish a repeatable means of approaching the development of vibrotactile belts. While we discuss most of the design issues in the context of developing vibrotactile belts, we are confident that these guidelines can be immediately extended to any wearable vibrotactile display technology.

10.4 Implementation

10.4.1 Form Factor

A belt's form factor ultimately determines its wearability and portability. To this end, we attempted to make the belt as robust and wearable as possible (see Fig. 10.4.1). The control box offers complete belt control along with wireless connectivity and battery power supply, and measures 8 cm by 4 cm by 2 cm. The individual tactor modules enclose a separate controller and a vibration motor, and measure 5.4 cm by 3.49 cm by 1.47 cm. The belt was designed to be lightweight (harness: 92.14 g; each tactor: 21.26 g; and controller: 95.68 g), comfortable and physically discreet.

The belt harness (flat nylon webbing) is easily adjustable to any waist size using plastic buckles while the tactors and control box are on pocket clips and can be adjusted appropriately per application, in seconds. This design was chosen over a Velcro based implementation (popularly encountered in our literature survey) to achieve better adaptability to different waist sizes; to hold tactors very close to the body during use; and to offer ro-



business and rigidity for real-world applications. The control box and the individual tactors are connected over a 4 -wire I2C bus that carries power along with the data and clock. This configuration, allows plug and play adding, removing and reconfiguring of tactors for scalability and reconfigurability.

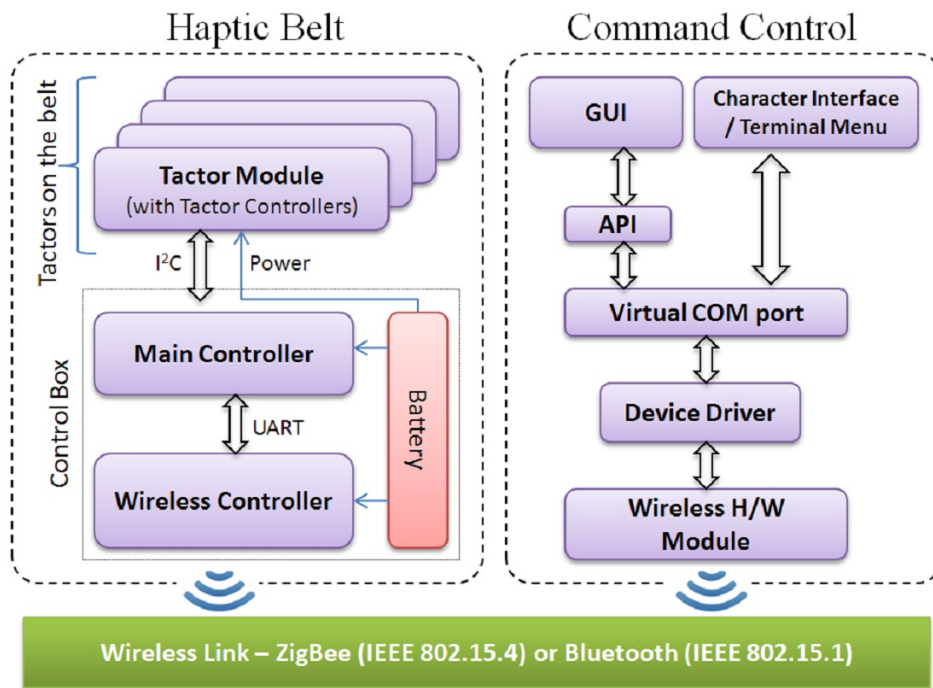
10.4.2 System Architecture

In order to provide two important functional requirements of expressiveness and scalability, we employ a network of distributed controllers. The hierarchical system level design of the belt, shown in Fig. 10.5, utilizes an independently functioning wireless main controller (Haptic Belt Controller) enclosed within the control box, and auxiliary controllers (Tactor Controllers) for monitoring and controlling each vibration motor, represented as tactors in Figure 10.5.

While the main controller offers connectivity to a command control center (PC or PDA), each tactor controller takes care of the micromanagement of vibrotactile cueing at each vibration motor. This multi level hardware processing buffers commands, and consequently allows for a higher performance and responsiveness of the system when compared to a centralized processing system. Each sub-system encapsulates its functionality locally

so that it provides functional independence from other sub-systems, all while achieving this with minimal data transmissions. Any shared data is stored centrally on the main controller and is distributed on power-up or redistributed after a configuration change. One of the important design requirements of a haptic belt, reconfigurability, is the capability to configure the belt's parameters easily with the bare minimum software tools. To this end, the belt connects through character terminal interface with Hayes AT command like serial communication interface.

10.5 Hardware Design



10.5.1 Control Box

The control box receives all control messages transmitted from the command control center (PC or PDA) to the haptic belt. As shown in Fig. 10.5, the most important components of the control unit are as follows:

10.5.1.1 Main Controller

A specific implementation of the popular Arduino Open Source hardware platform (based on Atmel ATMeg168 microcontroller), called Funnel IO , was used for the main controller.

10.5.1.2 Bus Communication

One of the most important requirements for the design of the haptic belt was the need to reduce the number of wires connecting factors. It was this constraint which led to the use of individual controllers at each of the factors. Complementing this choice, I2C offered the least number of wires with reliability. Thus, a four wire bus implementation, 2 wires for power, one for data (SDA) and one for clock (SCL), was adopted. The implementation allows up to 16 factors on the belt simultaneously.

10.5.1.3 Power Supply

Much consideration was given to the possible use time of the belt when specifying and sizing the power supply technology. Considering the space constraints, Lithium-polymer chemistry provides the most charge density for its size, and so a single cell 3.7V 800 mAh battery that allows up to 6 hours of continuous operation was chosen.

10.5.1.4 Wireless Hardware

Our performance requirements for the wireless module included transmission range of a large room, and the inclusion of a separate microcontroller to manage transmission without impacting general controller function. Either of two integrated wireless modules (Digi's XBee ZigBee module and Roving Network's RN-41 Bluetooth module) were chosen to connect to the Funnel board through a dedicated UART providing the necessary wireless connectivity and control.

10.6 Tactor Modules

As shown in Fig. 2, the tactor modules individually contain a microcontroller that negotiates its role with the main controller through the I2C bus. An Atmel ATtiny88 microcontroller forms the core of the tactor module. The PWM unit on the microcontroller is used for amplitude control and temporal rhythm generation, as described in Software Design (Section 10.7), while running independently from the main controller. A MOSFET driver provides the necessary switching between the digital output and the motor actuations. Six GPIO pins of the ATtiny88 are configured to read a DIP switch setting that assigns each tactor module's bus address. This address is used by the main controller to dynamically assign the I2C bus address at startup. This eliminates the need to reprogram all tactor modules for different applications/uses, thereby providing plug-and-play functionality. Vibrations are actuated through use of a 12 mm coin-type shaftless vibration motor, which has a rotational speed of 150Hz and a nominal vibration of 0.9g. The motors were mounted such that the vibration axis is parallel to human skin causing a net lateral vibration along the skin.

10.7 Software Design

The software components of our proposed design contain two important aspects: the firmware, which is programmed on the microcontrollers, and the User Interface (UI) that allows the design of vibrotactile rhythm patterns and access to the operational modes of the haptic belt.

10.7.1 Firmware

As explained earlier, the proposed haptic belt system makes use of a distributed microcontroller network framework with a separate main controller and the tactor microcontrollers for increased functionality and reliability. Below we discuss the important aspects of the firmware for the two controllers.

10.7.1.1 Main Controller Firmware

The main controller provides communication between the command control (through wireless protocols) and the factors on the belt. The main controller's firmware can be categorized into 7 primary functional areas as shown in Fig. 10.2.

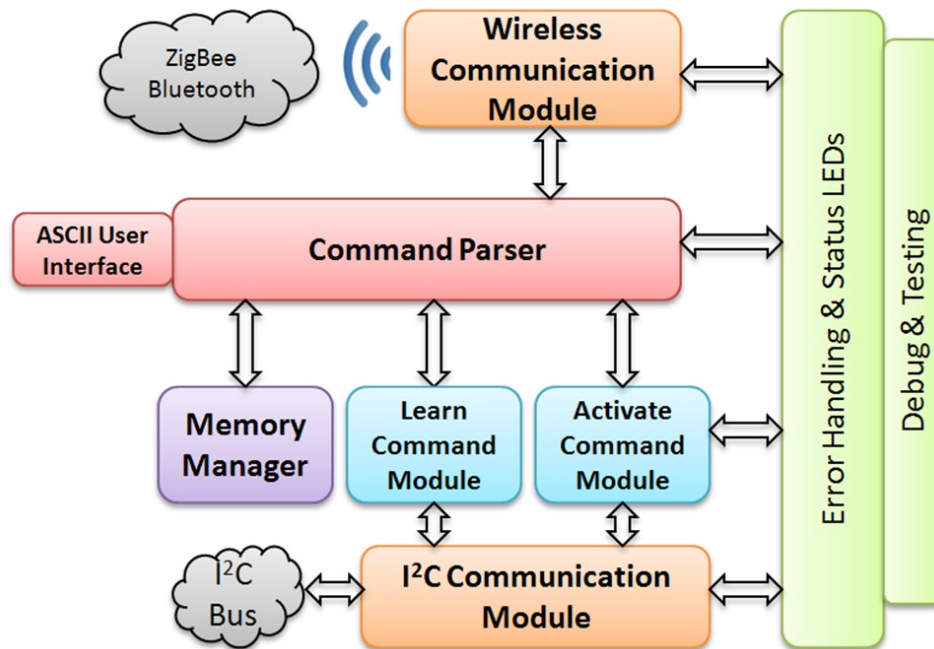


Figure 10.2: Main Controller implementation

- (a) *Wireless Communication Module*: All communication from the command control center (PC) is received through the ZigBee/Bluetooth wireless module. This module reads and writes data to and from the hardware buffers in a continuous loop. All data received is automatically sent to the Command Parser for further interpretation.
- (b) *Command Parser and ASCII User Interface*: This module provides four primary user modes, namely, a) new belt configuration; b) query current configurations; c) test vibrotactile patterns; and d) binary command mode. These modes allow the user to configure, use and debug individual factors and the belt as a whole.

- (c) *Learn Command Module*: In the proposed belt design, versatility is provided through user definitions of the Temporal Rhythm Unit (TRU) and Temporal Rhythm Sequence (TRS) (See Section VI.B). This module handles all the activities of the learning module while building rhythm pattern definitions (TRS and TRUs). This module also sends all new configurations to the Memory management module (via command parser) to be stored in the on-chip memory.
- (d) *I2C Communication Module*: This module is responsible for querying all tactors (or any devices) on the bus and stores their addresses into a data table. This module is also responsible for sending commands and receiving status codes from all the tactor modules.
- (e) *Memory Manager*: The ATmega168 controller has limited SRAM for runtime operations. The memory manager is implemented so that all rhythm definitions or text-based menus can be stored and retrieved from the on-chip Flash memory. The command parser handles the control flow to the memory manager.
- (f) *Activate Command Module*: This module handles the binary encoding of a tactor activate command. It packages the requested rhythm (TRS) and magnitude (TRU) with the appropriate cross-referenced tactor bus address and sends the command to the specific tactor for activation.

10.7.1.2 The Tactor Controller

The second part of the system includes the hardware for the tactor module. As shown in Figure 2, the tactor modules individually contain a microcontroller that negotiates its role with the main controller through the I2C bus. In this section, we provide details of the design choices and the implementation of the tactor.

The Vibration Motor: The vibrations are actuated in our tactor modules through the use of 12 mm coin-type shaftless vibration motors manufactured by Precision Microdrives, which have a rotational speed of 150Hz and a nominal vibration of 0.9g. These motors use an off-center mass to actuate vibrations. When the motor is powered, the rotation of the

shaft, and hence the eccentric mass, causes vibrations that is maximal along the direction perpendicular to the rotational axis. In our implementation, we mounted the motor in such a way that the vibration axis is parallel to human skin causing a net lateral vibration along the skin. Recently, we have incorporated cylindrical vibration motors that are mounted such that the vibration axis is perpendicular to the human skin. Experimental results with both have not revealed significant difference in their vibration conveyance.

Dimensions of the Vibration Signal The three important dimensions of the vibration signal, as delivered through our haptic belt, include a) the location of vibration, b) the amplitude of vibration, and c) the timing and temporal rhythm of vibration:

- (a) Location of Vibration: The location of vibration is reflected by the flexibility of our belt which allows the easy movement, addition or removal of tactors from the belt strap. This allows users of the belt to achieve any positional configuration that is desired for the application in focus.
- (b) Amplitude of Vibration: A pulse-width modulation (PWM) based amplitude control (similar to digital sound modulation) is incorporated to control the intensity of vibrations. The applied voltage is modulated in 20 microsecond intervals to achieve desired levels of intensity. Although the design allows for a much smaller resolution, human sensory mechanisms cannot practically distinguish them. Figure 10.3(a) and 10.3(b) shows the duty cycles under 2 different magnitudes, and Figure 10.3(c) and 10.3(d) shows sample waveform of the actuation signals generated by the tactor microcontroller under 25% and 75% intensities.
- (c) Timing and Temporal Rhythm of Vibration: Temporal rhythm patterns refer to the actuation of the vibrators as discrete time pulses. Readers should not confuse these temporal pulses with the pulse-widths used for magnitude control as described in the section above. The vibration of any tactor on the belt is divided into discrete temporal time events referred to as the Temporal Rhythm Unit (TRU), where each TRU is 50 ms long. The choice of 50 ms per TRU was determined experimentally, where we

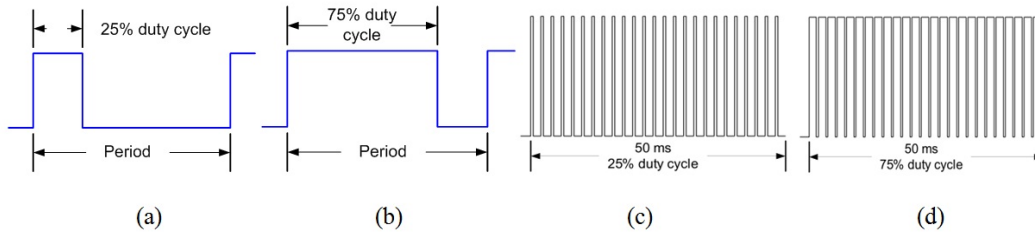


Figure 10.3: (a) 25%; (b) 75% Pulse-width modulation; (c) and (d) Vibration motor magnitudes of 25% and 75% achieved using duty cycles with 25 pulses over a 50 ms vibration period.

found that any vibration pulse of duration lesser than 50 ms was not perceivable to the participants. All vibrations are defined in terms of the number of TRUs for which the vibration motor should be ON or OFF. This sequence of ON-OFF patterns forms a Temporal Rhythm Sequence (TRS). Note that any ON TRU can be controlled in magnitude by incorporating the PWM amplitude control as explained in the previous section. Two sample TRS are shown in Figure 10.4 below. The first sequence has 5 TRUs, ON-OFF-ON-ON-OFF, totaling 250ms with amplitude of 100%. The second TRS has 4 TRUs, ON-OFF-ON-OFF, totaling 200ms with amplitude of 50%.

The tactor controller firmware communicates directly with the main controller firmware as a slave device over the I2C bus and maintains the PWM timing for the local vibration motor. A two-byte command structure is used between the main controller and the tactors. Similar to the main controller firmware, the functionality of the tactor controller can be categorized into five important roles (Fig. 10.5). While the communication module and command parser are similar as above, the memory module and the low-level hardware module (PWM module) form the critical components of the tactor module. The memory manager module is responsible for temporarily storing the definition of the TRU and TRS that are sent over to the tactors at boot up. At run time, a two-byte activate command selects the appropriate TRU and TRS for each tactor, which the PWM module executes on the vibration motor.

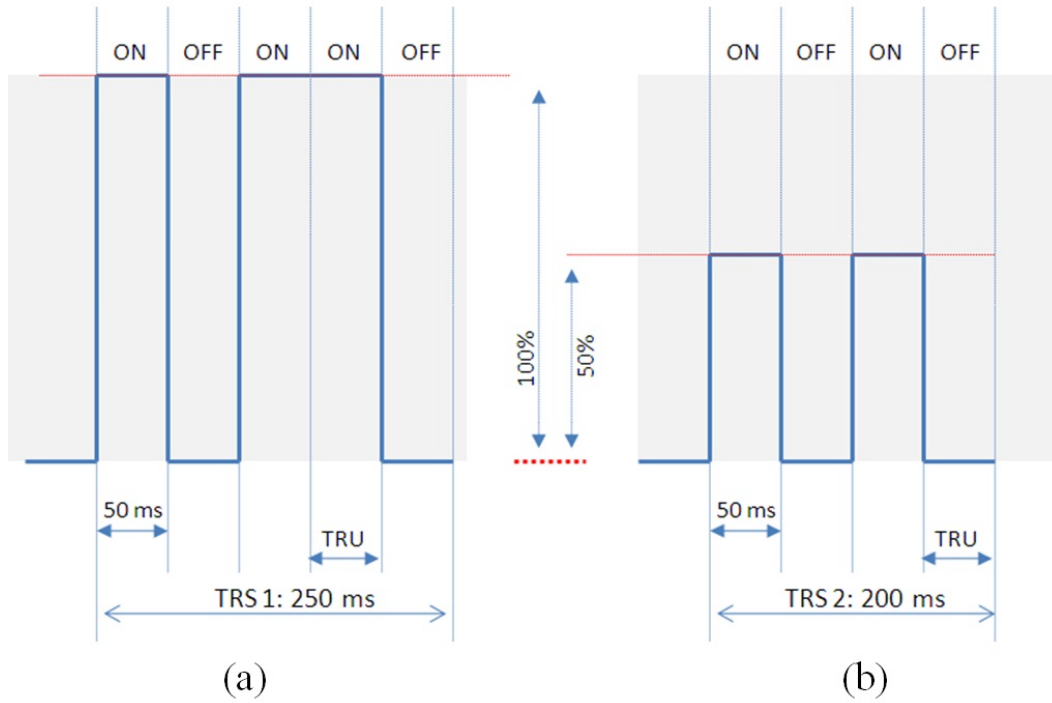


Figure 10.4: Sample Temporal Rhythm Sequences (TRS) with different magnitudes of vibration encoded on the Temporal Rhythm Units (TRU) (a) 100% Magnitude, (b) 50% Magnitude.

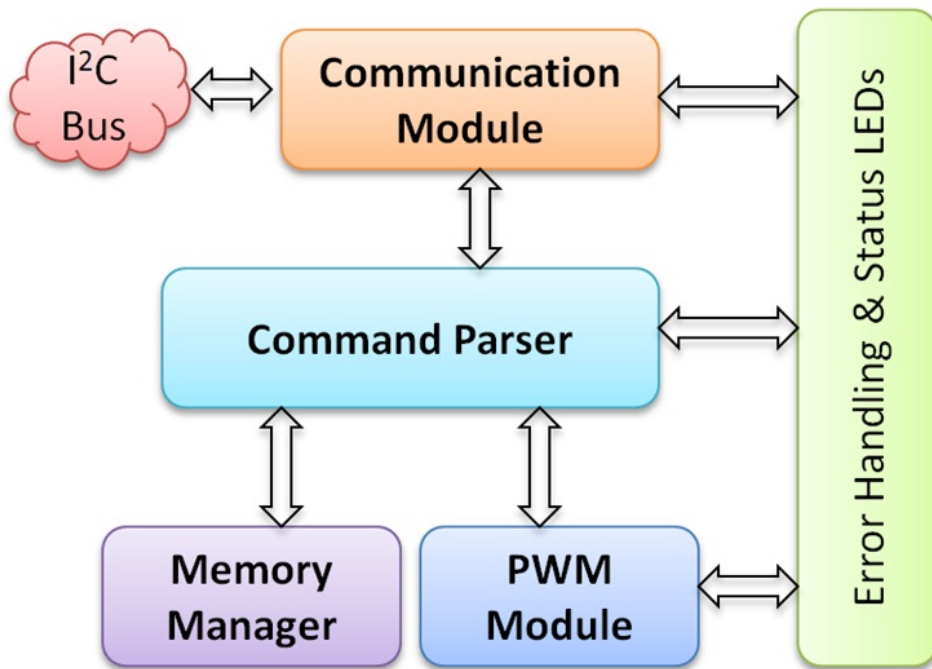


Figure 10.5: Tactor Controller implementation

10.8 User Interface

The user interface on the haptic belt currently supports two complementary formats: 1) A console based Hayes AT command set-like interface for quick access to all functionalities of the belt, and 2) an Application Programming Interface (API) for more advanced programming in higher level languages and for Graphical User Interface development.

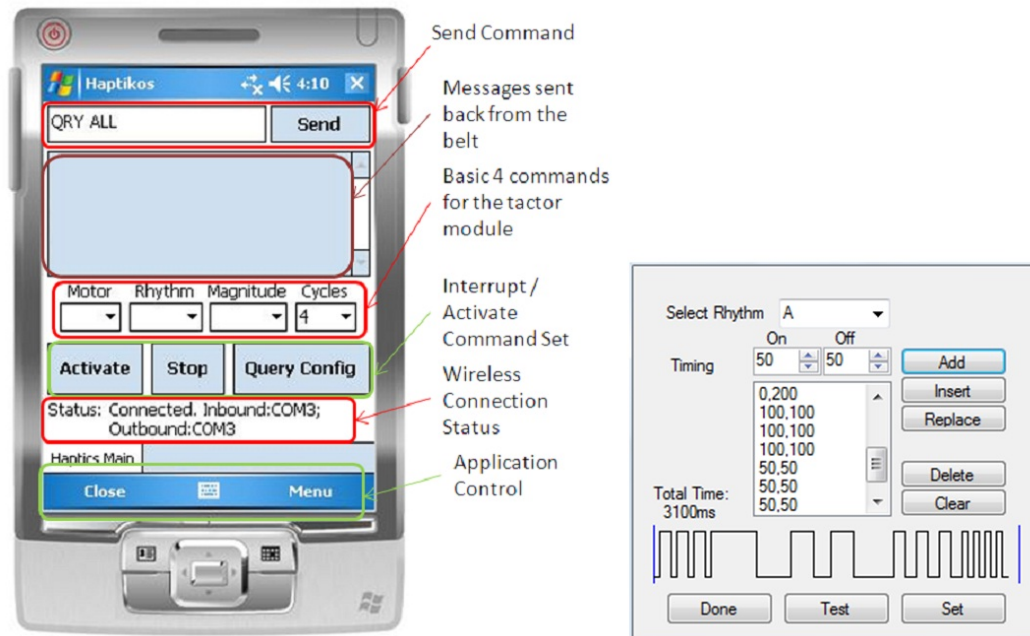


Figure 10.6: a) Graphical User Interface on a Portable Platform. b) Temporal Rhythm Sequence (TRS) Design Interface.

Currently, we have implemented a PC based and a PDA based GUI for controlling and configuring the haptic belt. The GUI allows the design of complex vibrotactile rhythm and spatio-temporal patterns. The API has the portability of supporting a wide range of ubiquitous computing platforms, including mobile devices. In Fig. 10.6(a), we show the PDA interface with highlights on some of the important features. While the functionalities on the PC are very similar to the PDA version, additional features on the PC allow easy configurability of the belt. The example feature shown in Fig. 10.6(b) shows the setup used for designing a TRS. Users can select specific rhythms and vary the TRUs appropriately based on the application. Each TRU is 50 ms long, and the entire TRS can be a maximum

of 3 seconds long. The interface allows a user to compose a pattern of patterns by interleaving TRSs. The case study discussed later in this paper used these controls to design haptic patterns related to the study. Our graphical interface is similar to other haptic pattern authoring software like posVibEditor [359], and provides a framework for rendering digitally modulated vibration patterns. As there is a lack of standardization and open sourcing among authoring tools for vibrotactile patterns, our efforts are unique in that our work is available for download, as of this publication.

10.9 Experiments

The experiments presented here tested the haptic belt system for its use in conveying non-verbal cues, specifically cues pertaining to where communicators are located in front of the user in terms of direction and distance. Figure 10.1 shows seven vibrators located in the front of the user in the form of a semicircle. This setup allowed us to focus on accurately assessing the capabilities of the haptic belt.

10.9.1 Experiment 1: Localization of Vibrotactile Cues

Prior work [352] showed that reasonable localization accuracy-between 80% to 100% accuracy depending upon factor location-was possible with a belt design similar to what we presented above. Our experiment is similar, but offers a few variations to verify the results obtained in [352]. Subjects: 10 subjects (8 males and 2 females), of ages between 24 and 59, participated in this experiment. One of the subjects was blind; the rest were sighted. Subjects had no known deficits related to their tactile sense of the waist area. Further, no subjects had prior experience with haptic belts, but all subjects had some exposure to vibrotactile cues (e.g., vibrations of a cell phone).

10.9.1.1 Apparatus:

The haptic belt described earlier was used for this experiment. Vibratory signals were 600 ms in length, and had a frequency and intensity well within the range of human perception. In contrast to [352], cues are longer-600 ms compared to 200 ms-and we do not use head-

phones to mask subtle vibration noise, nor do we randomly vary intensity with each cue; the reason for these changes is that we are mostly concerned with how the belt as a complete system accomplishes non-verbal communication, rather than the spatial acuity of the waist. Hence, if a specific intensity of vibration feels different around the waist, and some vibrations can be heard, and if these cues help in factor localization, then this redundant information should only add to the usability of the system.

10.9.1.2 Procedure:

Subjects put on the haptic belt over their shirt and around their waist such that the middle tactor (#4) was centered at their navel, and the endpoint tactors (#1 and #7) were at their left and right sides, respectively. As the belt has LEDs that light up to indicate tactor activation (used for testing the belt), subjects were instructed to not look down at the belt any time during the experiment. Next, subjects were familiarized with tactor numbering: the experimenter activated tactors in order from #1 to #7, and spoke aloud the number of the activated tactor. This process was repeated twice for each subject.

The training phase involved 35 trials where each tactor was randomly activated 5 times (with approximately 5 seconds between tactor activations) and subjects had to identify the number of each activated tactor. A visual guide was provided for subjects to help recall tactor numbers; this guide was a white board with a drawing of a semicircle (the belt) and the numbers 1 through 7 (tactors) on the belt. Feedback was given during the training phase to correct wrong guesses. The testing phase was similar to the training phase, but involved 70 trials where each tactor was randomly activated 10 times, and feedback was not provided. Subjects stood during the entire experiment.

10.9.1.3 Results:

The localization accuracy for each tactor (number of times identified correctly out of the total number of times activated) was averaged across subjects and is shown in Figure 10.7 (indicated by the dots centered within each error bar), where error bars indicate 95% con-

fidence intervals. The overall localization accuracy across factors and subjects was (92.1 7.0)%.

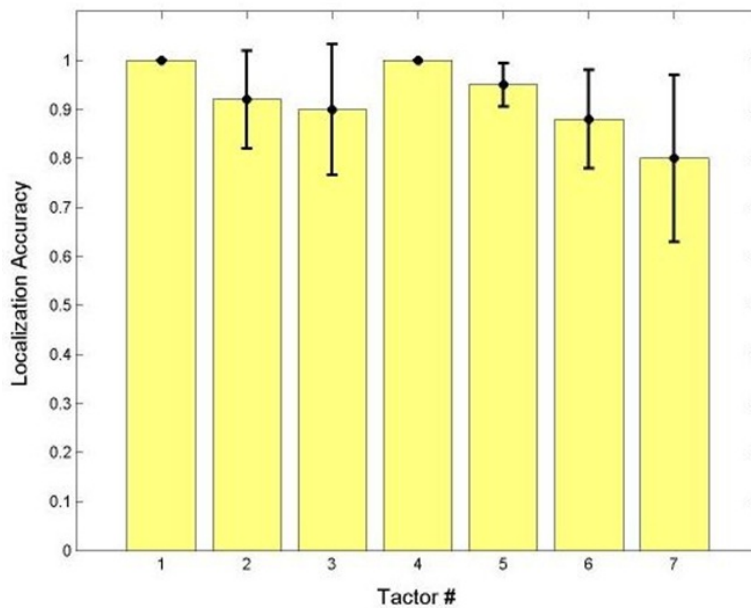


Figure 10.7: Experiment 1 Results: Mean Localization Accuracy for each Tactor, Averaged across Subjects, with 95% Confidence Intervals

10.9.1.4 Discussion:

An overall localization accuracy of (92.1 7.0)% (an improvement over that of [352]) is promising and shows that our prototype haptic belt can be reliably used to indicate the direction of someone in the user's visual field. Moreover, 100% of misclassifications were off by a single tactor location; hence, even when users made a mistake in localizing an activated tactor, they still had a very good idea of the general direction of someone in their visual field.

We hypothesize that the increase in accuracy is largely due to greater cue duration (600 ms as opposed to the 200 ms used in [352]); it is well known that larger cue durations make localization easier [353]. Moreover, redundant information provided by the belt, such as subtle audible cues when tactors are activated, could have helped as well. Subjects found tactors closer to the midline easier to localize, which agrees with the results found in the

literature where spatial acuity improves near the sagittal plane [352] [353] given that spatial acuity is better at anatomical reference points-in this case, the navel.

It is hypothesized in [352] that the tactors at the end of the semicircle, which rest at the sides of the torso, act as landmarks and are easier to localize; but in our experiments, we noticed that tactor #1 could be localized more accurately than tactor #7, as shown in Figure 4. We are investigating this asymmetric result.

10.9.2 Experiment 2: Signal Duration as Cue for Distance

Experiment 2 included two sub-experiments to examine the use of vibrotactile cues to indicate both direction and distance. Experiment 2A focused on how well subjects could perceive cue duration, regardless of tactor location. Experiment 2B tested how well subjects could perceive both tactor location and cue duration at the same time.

10.9.2.1 Subjects:

The ten subjects introduced for Experiment 1 also performed Experiment 2A and 2B.

10.9.2.2 Apparatus:

The belt and signal properties were identical to those of Experiment 1 with the exception of signal duration. For Experiment 2A and 2B, signal durations of 200 ms, 400 ms, 600 ms, 800 ms and 1000 ms were used. These durations may refer to any distance in the implementation of the system; e.g., less than 2 ft (1000 ms), 2 ft to 4ft (800 ms), 4 ft to 6 ft (600 ms), and so on.

10.9.2.3 Procedure:

In the first part of Experiment 2A, subjects were familiarized with the five cue durations. All five durations were delivered to the user at each of the seven tactors from #1 to #7, in order. The training phase for Experiment 2A involved 35 trials where each tactor was randomly activated 5 times (one time for each of the 5 durations) with approximately 5 seconds between tactor activations. Subjects were instructed to guess only cue duration.

The testing phase involved 70 trials with each factor activated twice for each duration. As in the training phase, subjects had to guess the duration of the cue, but no feedback was provided. Immediately following Experiment 2A, subjects began Experiment 2B. First, the familiarization and training phase of Experiment 1 were repeated. The testing phase involved 70 trials similar to 2A, but subjects now had to guess both cue duration and factor location. As in Experiment 1, subjects stood the entire experiment, had access to a visual guide and were told not to look at the belt.

10.9.2.4 Results:

Classification accuracy of duration (number of times identified correctly out of the total number of times used) was averaged across subjects and is shown in Figure 10.8 (indicated by the dots centered within each error bar), where error bars indicate 95% confidence intervals. Note that the x-axis of Figure 5 lists durations as #1 (200 ms), #2 (400 ms), #3 (600 ms), #4 (800 ms) and #5 (1000 ms). The results for both Experiment 2A and 2B are included in Figure 5. The overall classification accuracy of duration across factors and subjects was (73.3.6)% and (67.11.8)% for Experiment 2A and 2B, respectively. There were not any noticeable differences in classification accuracy of duration between different factor locations in either part of the experiment.

10.9.2.5 Discussion:

In Experiment 2, subjects were able to easily identify durations of 200 ms and 400 ms, most likely due to their short length. However, subjects had difficulty distinguishing between 600 ms, 800 ms and 1000 ms. Two subjects suggested that a logarithmic scale of 200 ms, 400 ms, 800 ms, 1600 ms, and so on, might improve recognition. However, longer cues slow down use of the system, making it more difficult to use in real time. Another option would be to use fewer cues (e.g., 200 ms, 500 ms and 1000 ms) to provide only coarse distance information. Regardless, the overall classification accuracy of duration at (73.3.6)% is impressive, and accuracies for longer durations are satisfactory. The skill of subjects at classifying lengths of vibrations varied, resulting in large variations in classification ac-

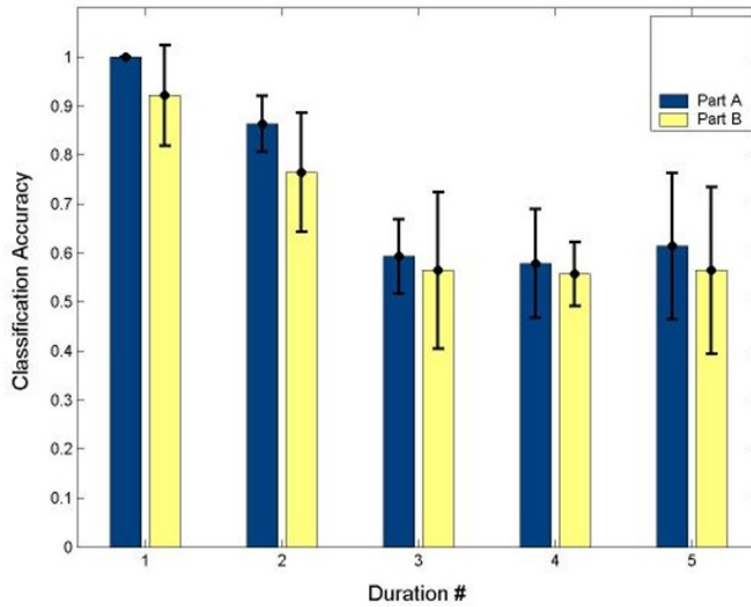


Figure 10.8: Mean Classification Accuracy of Duration, Averaged across Subjects and Tactors, with 95% Confidence Intervals. Durations listed in figure correspond to 200 ms (#1), 400 ms (#2), 600 ms (#3), 800 ms (#4) and 1000 ms (#5)

accuracy for longer cue durations (see Figure 5). In any case, 94.7% of misclassifications were off by only 200 ms (5.3% of misclassifications were off by 400 ms), which shows that subjects were quite accurate with their estimates.

In Experiment 2B, overall classification accuracy of duration dropped to (67.1 ± 11.8)%. We hypothesize that this small drop in mean accuracy, as well as an increase in variance, was due to the cognitive load of having to attend to both vibration duration and tactor location. In any case, overall accuracy is still satisfactory, and 89% of misclassifications were off by only 200 ms (11% of misclassifications were off by 400 ms). Overall tactor localization accuracy for Experiment 2B was (94.3 ± 5.7)% (averaged across subjects, tactors and durations), which is similar to the localization accuracy of (92.1 ± 7.0)% found in Experiment 1. Once again, 100% of misclassifications were off by a single tactor location. We conclude that tactor locations are still easy to perceive even when cue length varies and attention must be divided between cue duration and location.

10.9.3 Experiment 3: Vibrotactile Rhythm as Cue for Distance

As an alternative to vibration duration, we also explored the use of vibrotactile rhythm to deliver distance information.

10.9.3.1 Tactile Rhythm Design

The tactile rhythms used in our experiments were motivated by results reported in [344], where just noticeable differences of vibrotactile duration were assessed. Subjects perceived pulses of duration below 100ms as a poke or nudge. Between 100ms to 2000ms, the just noticeable difference is an increasing curvilinear function of duration; although between 100ms to 500ms, the function is approximately linear. Based on these results, Geldard [344] recommended three durations, specifically 100ms, 300ms and 500ms, for accurate identification by subjects.

Tactile rhythms delivered using a vibrotactile belt were used in [348] to convey distance information during waypoint navigation. Time between vibratory pulses was varied using one of two schemes: monotonic (rate is inversely proportional to distance) or three-phase-model (three distinct rhythms mapped to three distances). Distinct tactile rhythms are promising for use with multidimensional tactons [360] [361], which are vibratory signals used to communicate abstract messages [361] by changing the dimensions of the signal including frequency, amplitude, location, rhythm, etc. Based on pilot test results, we chose to pursue distinct rhythms over monotonic rhythms as users find it difficult to identify interpersonal distances using monotonic rhythms as the vibratory signal varies smoothly with changes in distance.

We conducted pilot studies to determine rhythm patterns that are convenient for users to identify vibratory rhythms. Through use of a vibrotactile belt, we evaluated use of five rhythms, each 10 seconds in length: 50ms vibrotactile pulses separated by pauses of length 50ms, 100ms, 300ms, 500ms and 1000ms. Subjects found rhythms with pauses of 100ms, 300ms and 500ms difficult to discriminate between. Based on these findings, we

selected the four rhythms depicted in figure 10.9; this design includes more separation of pauses within 100ms to 500ms, and a small increase of 1000ms to 1200ms (much longer durations may be too time consuming for communication [344]). In the Social Interaction Assistant, these four tactile rhythms are mapped to interpersonal distances corresponding to intimate, personal (close phase), personal (far phase) and social (close phase) space respectively.

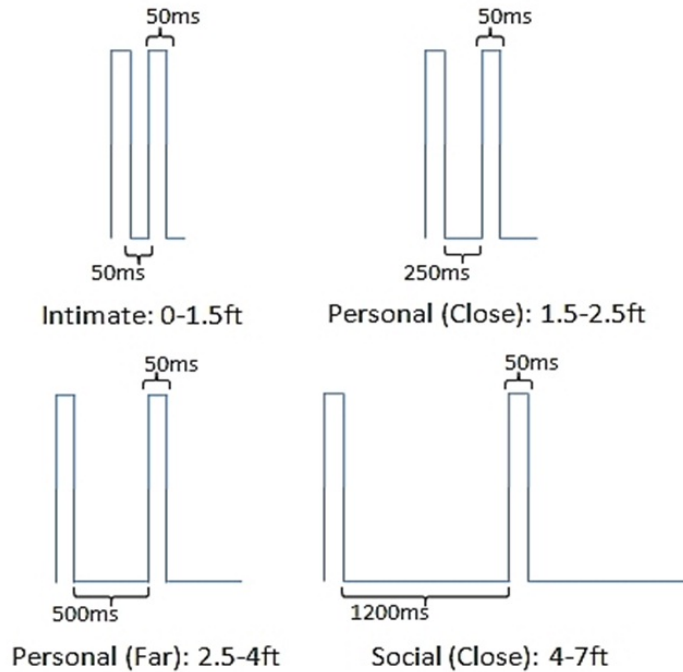


Figure 10.9: The on/off timing values of the four tactile rhythm designs, and corresponding distances, used in the experiment.

10.9.3.2 Experiment

Aim: The aim of this experiment is to evaluate participants' performance identifying the tactile rhythms of figure 1 as they relate to interpersonal distances. Moreover, to ensure that the proposed tactile rhythms do not hamper subjects' ability to localize vibrations, as evaluated in previous work [252] to convey directions, we evaluate how well subjects can identify both cues as conveyed through tactons.

Hypotheses: (1) Subjects will achieve at least 90% accuracy at identification of tac-

tile rhythms; (2) Subjects will achieve at least 90% accuracy at identification of vibration locations; (3) Subjects will achieve at least 80% accuracy at identification of complete tactons; (4) Subjects' ability to localize vibrations will depend on the location of the vibration motor (tactor) around the waist; (5) Subjects' ability to identify tactile rhythms will depend on the type of rhythm; and (6) Subjects' ability to localize vibrations will not depend on rhythm type, and vice versa.

Subjects: 11 males and 4 females of ages 22 to 60 (avg. 32) participated; one subject is visually impaired.

Apparatus: An elastic vibrotactile belt [252] was used for this experiment. The design of the belt was based on the experiments of Cholewiak, et al. [355]. The belt consists of 7 tactors equidistantly placed in a semi-circle with the first, fourth and seventh tactor at the user's left side, navel, and right side, respectively. Each tactor consists of a pancake motor of diameter 10mm and length of 3.4mm, and operates at 170Hz.

Procedure: Subjects wore the belt underneath their clothing and sat during the entire experiment. Subjects had access to visual guides—a semi-circle with tactors #1-7 drawn and interpersonal distances labeled as rhythms #1-4—to recall tactor and rhythm numbers, respectively. First, subjects were familiarized with vibration location as it pertains to direction. Each tactor was vibrated for 3 seconds, and the tactor number was called out by the experimenter. Next, subjects were familiarized with tactile rhythms. Each rhythm was presented for 7 seconds through the fourth tactor at the navel, and the rhythm number was called out by the experimenter. Next, subjects began the training phase where they were asked to identify the direction (through the location of the activated tactor) and distance (through the type of rhythm) indicated by each tacton. All 28 tactons (4 tactile rhythms at 7 different locations/tactors) were randomly presented for 10 seconds each. Subjects were encouraged to respond before the 10 seconds ended. Subjects had to achieve a recognition accuracy of 80% or more on each tacton dimension to proceed immediately to the testing phase; otherwise, the training phase was repeated (only 6 subjects had to repeat training, and all passed on the second try). The experimenter corrected wrong guesses and confirmed

correct guesses. The testing phase was similar to the training phase, except no feedback was provided by the experimenter concerning right or wrong guesses, and each of the 28 tactons was randomly presented 3 times for a total of 84 trials.

Results: The overall recognition accuracy follows (See Figure 10.10): location (mean: 95%, SD: 4%), rhythm (mean: 91.7%, SD: 5.7%) and both (mean: 87%, SD: 8.5%). These results support hypotheses (1)-(3), and show that, overall, subjects had little difficulty in recognizing rhythms and locations as they pertain to distance and direction, respectively. Feedback from participants after the experiment further supported this. From herein, reported ANOVA results are from a two-way ANOVA on complete tacton recognition accuracy through location and rhythm. The overall recognition accuracy of each tacton location is shown in figure 2. Subjects felt that the vibrations of tacton #1 (left side), #4 (navel) and #7 (right side), were easier to localize compared to tacton #2, #3, #5 and #6. This result is easy to explain as spatial acuity is better at anatomical reference points [355]. Although figure 10.10 does show a very small difference between recognition accuracies, which supports what subjects reported, there was no significant difference between recognition accuracy of tacton locations [$F(6,1232)=1.96, p=0.068$], hence hypothesis (4) cannot be accepted. The overall recognition accuracy of rhythms is shown in figure 10.11. Subjects felt that rhythm #2 (personal-close) and #3 (personal-far) were more difficult to identify than rhythm #1 (intimate) or #4 (social-close), which is supported by figure 3. A significant difference between recognition accuracy of rhythms [$F(3,1232) =5.70, p=0.001$] supported hypothesis (5). No interaction was found between location and rhythm for recognition accuracy of complete tactons [$F(18,1232)=0.91, p=0.569$], supporting hypothesis (6).

After the experiment, subjects filled out 10-level Likert scales-1 (lowest) to 10 (highest). Subjects rated their ability to localize vibrations (mean: 8.4), identify rhythms (mean: 7.4), intuitiveness of location to convey direction (mean: 9.7) and intuitiveness of rhythm to convey distance (mean: 8.9). Overall, subjects felt that they could accurately identify the proposed tactons, although identifying direction was easier than distance, and both schemes were intuitive.

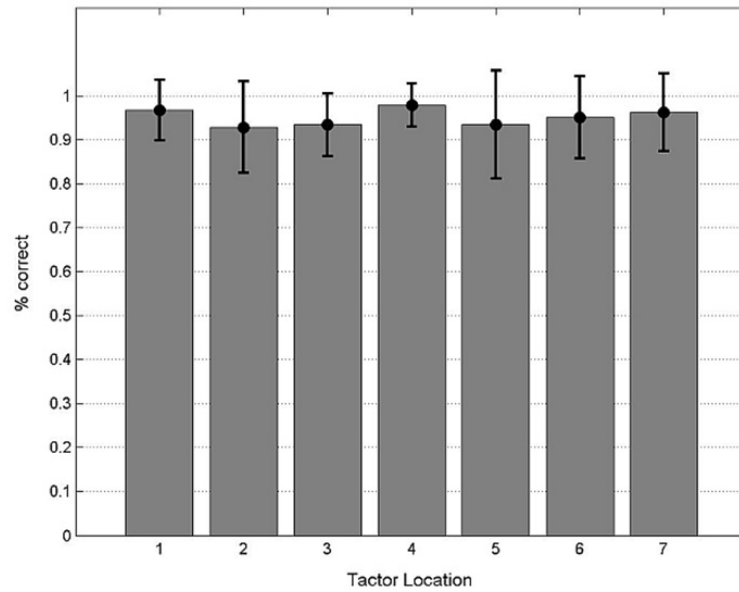


Figure 10.10: Overall direction recognition accuracy of each factor location with standard deviations.

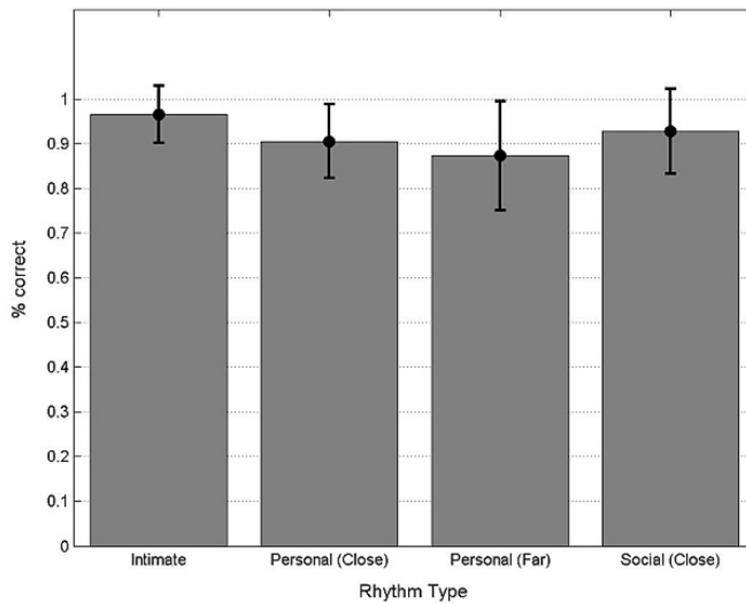


Figure 10.11: Overall distance recognition accuracy of each rhythm type with standard deviations.

10.10 Other Applications for the proposed technology

As briefly mention in Section I, haptic belts provide for a wide array of useful application areas. Here, we elaborate on three specific application areas of haptic belts.

10.10.1 Navigation and Spatial Orientation

The most prevalent application for haptic belts is in the development of navigational [350] [348] [349] and spatial orientation aids [362] [356], which have been well explored in academic research and hobby development. Conceptually, a haptic belt designed for navigation will make use of a positioning system, either absolute (e.g., GPS, GLOSNASS or Galileo) or relative (e.g., Inertial Navigational Units), along with a map of the locality to guide the user from their current location to a desired destination as shown in Figure 5. Since the vibrotactile actuators are mounted around the waist of the user, directional information is conveyed through activation of the appropriate motor.

Humans generally work or move about by using geographical references, and without them, it is easy to become disoriented. Haptic belt solutions can be used for these applications, along with absolute positioning sensors, to determine and convey specific reference planes; e.g., a gyro based artificial horizon in the case of pilots [353], and the direction of Earth in the case of astronauts [354].

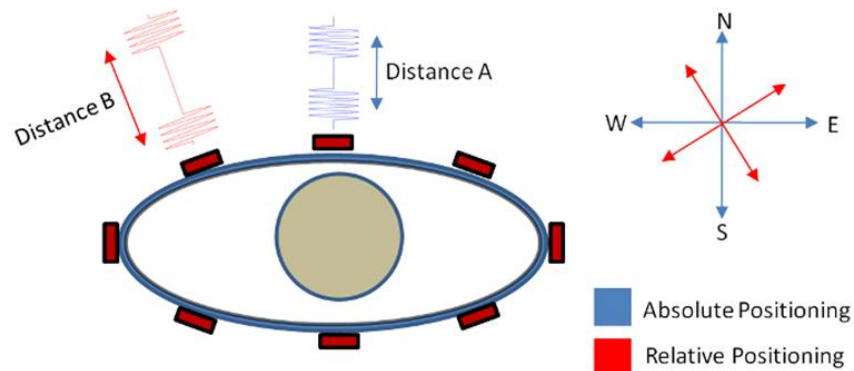


Figure 10.12: Application of haptic belt as a navigational aid.

10.10.2 Interpersonal Social Communication

At the intersection of research areas in social interaction, and that of human machine interfaces, lies the relatively unexplored field of interpersonal social communication technologies [82]. Affective Computing [363] has been exploring various sensor and actuator technology solutions toward better communication of social interpersonal signals with non-verbal communication cues [189]. With the human body having the largest sensory and perceptive surface, haptic displays, like our proposed vibrotactile belt, provide an opportunity toward enabling such communication of social signals. Over the last year, we've been working toward a wearable system capable of communicating interpersonal positions of interaction partners to someone who is visually impaired. In [252] [364], we developed an assistive technology that is capable of determining, and conveying through a haptic belt, the interpersonal distance and orientation of an individual who is standing in front of a user who is blind. Computer vision techniques are used to extract the interpersonal distance and orientation with respect to the user. The relative direction and distance of an interaction partner is conveyed through the location and rhythm, respectively, of a vibration around the user's waist.

10.10.3 Generic Information Communication

Given proper functionality, specifically expressiveness, haptic belts can be used to convey generic information through tactile icons [341], or tactons. Tactons are abstract, tactile messages, where meaning is mapped to the dimensions of the vibrations. Research has shown users' ability to quickly learn and interpret tactons [361]. For example, [360] demonstrated that anesthesiologists have enhanced situational awareness toward adverse medical events when a haptic belt is used to monitor a patient's physiological data in the operating room. Similarly, a haptic belt can be used to reduce visual and auditory cognitive load in military personnel [362], especially pilots and dismounted soldiers within a combat zone. A haptic belt could be used to communicate battlefield situational awareness when combined with military intelligence information and radio technology. Likewise, situational awareness,

spatial orientation and navigational information can be combined to provide real time location information through a haptic belt to emergency responders who have entered a low visibility and hazardous environment where visual and auditory modalities are overloaded with information.

10.10.4 Case Study: Waist-worn Vibrotactile Display for Pedagogical Application for Choreographed Dance

To evaluate the vibrotactile belt's three important design parameters of usability, functionality and performance, we conducted a case study in which the belt was used for a novel pedagogical application under realistic conditions. A two-fold, quantitative and subjective analysis was conducted to evaluate real-world usability issues. In [365], we demonstrated the general usability of the belt through a pilot study from the user's perspective, but the functional and performance metrics were not evaluated. In this paper, we delve into the details of the belt's evaluation through its use as a research platform for a novel application of teaching choreographed dance. While the usability analysis was done from the user's perspective, the functional and performance analyses were done by an independent researcher who designed and executed the dance study. The researcher was not part of the development team and evaluated the proposed belt as a research platform to impart choreography of simple dance steps to a mixed group of participants with and without dance experience. It is important to note that the researcher who worked with the proposed belt (a) had never used the vibrotactile belt; (b) had limited experience with haptic devices, and had never designed vibrotactile spatio-temporal patterns; and (c) had a novel application design with specific research objectives of determining how effectively choreography can be achieved through wearable vibrotactile devices.

10.10.4.1 Related Work in the Use of Vibrotactile Cues for Teaching Dance

In the literature, vibrotactile stimulation to elicit motor movement can be divided into two approaches: feedback-based and instruction-based, both of which are relatively new and unexplored. While feedback based approaches track human body motion and provide feed-

back whenever there is a deviation from a predefined path [366] [367], instruction-based methods assign specific body movements to predesigned vibrotactile patterns, and expect subjects to memorize this mapping. In [368], Drobny et al. developed a wireless sensory system placed in the shoes of ballroom dancers. By measuring the force of taps, the system recognizes any missteps and emphasizes beats acoustically to help partners get back in sync. While this study was centered on a feedback based learning system, this case study uses an instruction-based method for teaching dance (similar to various other pedagogical applications targeting physical activities such as snowboarding [369], bowling [370] [371] and swimming [372]), where predefined spatio-temporal vibration patterns require participants to demonstrate specific movements. To the best of the authors' knowledge, the only other work that explores instruction-based vibrotactile cues for teaching dancing is an approach by Nakamura et al. [373] where vibrotactile cues instructed arm movements for traditional Japanese folk dance. Unfortunately, the paper does not describe any of the proposed vibrotactile cues, and no statistical analysis was presented. Note that the cues proposed here are for basic dance movements only; more complex dance movements will require further exploration by dance experts on the possible redesign of vibrotactile stimulators to be placed in strategic locations on the body.

10.10.4.2 Subjects

11 males and 2 females of ages 21 to 60 (average: 30) participated in the dance study. No subjects had any tactile impairment around their waist. 5 subjects had never danced before, 4 subjects had less than one year of dance experience, and 4 subjects had a least 5 years of dance experience. The dance participants provided data for analyzing the usability of the belt, while the independent researcher offered evaluations for the functionality and performance of the belt. Although we would have preferred several researchers and/or developers to assess the usability of our belt through their own novel applications and user studies, this was not feasible due to time limitations and the need for a specific target application.

10.10.4.3 Procedure

The belt was configured with 8 tactors placed equidistantly around each participant's waist. Fig. 10.13 shows the configuration with tactor #1 is at the user's left side, tactor #3 at the user's navel, tactor #5 at the user's right side and tactor #7 at the user's spine.

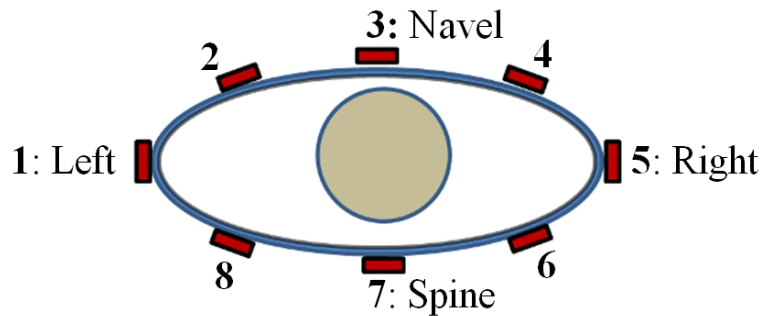


Figure 10.13: Arrangement of 8 Tactors around the Waist.

Subjects were informed that the purpose of the experiment was to assess how well they can recognize vibrotactile cues. They were not told that they would be learning basic dance moves to avoid giving any advantage to those with prior dance experience. Subjects were given instructions regarding how to put on the belt, and were told to move tactors along the length of the belt to match the configuration shown to them on a printed paper (same as Fig. 10.13). First, subjects were familiarized with the different vibrotactile patterns and the corresponding body movements (see Table II). Next, participants began the training phase where they were asked to feel a vibrotactile pattern and perform the associated movement, then return to the starting position. 24 trials (12 vibrotactile patterns each presented twice) were randomly presented. Subjects were encouraged to respond within ten seconds. Subjects were required to achieve above 70% accuracy in order to move on to the testing phase. The testing phase consisted of two parts. In the first part, the testing phase was similar to training phase, but with 48 trials and no feedback. Before the second part of the testing phase, participants performed another familiarization phase to help them learn how to link individual moves. In this familiarization phase, participants performed 11 moves in sequence. Finally, participants performed two different dance sequences: a

modified box step and a modified electric slide. The modified box step was repeated once, and consisted of the following vibrotactile patterns of Table II, in order: A, B, J, I, F, E, K, and L, as shown in Fig. 10.14. The modified electric slide was not repeated, and consisted of the following patterns of Table II, in order: J, I, J, I, K, L, K, L, F, E, A, B, B, A, E, F, K, L, K, L, J, I, J, I, A, B, F, E, F, E, A, and B, as shown in Fig. 8. A pause of 2 seconds was given between the pattern presentations. During this phase, no feedback was given to participants regarding right or wrong movements.

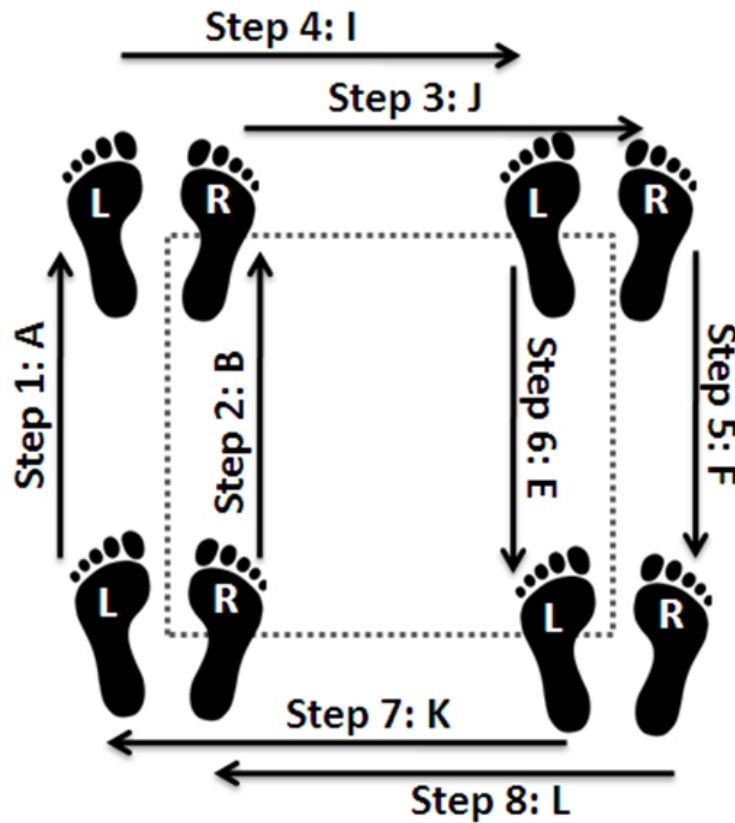


Figure 10.14: Modified Box Dance.

The independent researcher was given oral instructions (30 minutes) on the components of the belt and its complete operation including its software. The researcher was then allowed to configure the belt to his application. It took him 20 minutes to affix eight tactors and a control box, wire them together, and place them in the desired configuration. It took about 15 minutes to implement the vibrotactile cues for basic dance movements (left

Table 10.2: Foot Steps Involved in the Choreographed Dance Movements.

ID	Movement	Vibrotactile Pattern
A	Left foot forward (small step)	1 - 2 - 3
B	Right foot forward (small step)	5 - 4 - 3
C	Left foot forward (long step)	7 - 8 - 1 - 2 - 3
D	Right foot forward (long step)	7 - 6 - 5 - 4 - 3
E	Left foot back (small step)	1 - 8 - 7
F	Right foot back (small step)	5 - 6 - 7
G	Left foot back (long step)	3 - 2 - 1 - 8 - 7
H	Right foot back (long step)	3 - 4 - 5 - 6 - 7
I	Left foot right	1 - 2 - 3 - 4 - 5
J	Right foot right	3 - 4 - 5
K	Left foot left	3 - 2 - 1
L	Right foot left	5 - 4 - 3 - 2 - 1

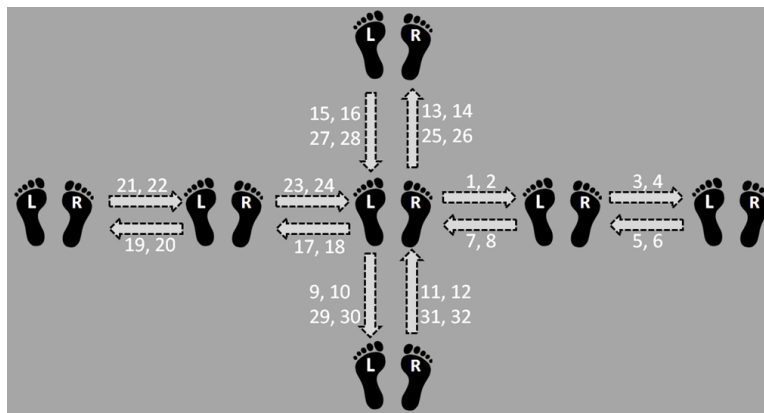


Figure 10.15: Modified Electric Slide Dance.

foot forward, right foot forward, etc.), and another 20 minutes to concatenate them into two dance sequences (modified box step and modified electric slide).

10.10.4.4 Aim

To evaluate the usability of our vibrotactile belt, we relied upon survey questionnaires completed by participants after the experiment. To evaluate the functionality and the performance, we relied upon the survey questionnaire completed by the researcher who worked with our team on conducting the experiment. Both surveys asked questions that directly or indirectly captured the various elements provided in Table 1. The usability questions for participants were:

- Q1. How easy was it to put on the belt and adjust the location of the vibration motors?
- Q2. How easy was it to take off the belt?
- Q3. How easy was it to recognize vibrotactile patterns corresponding to specific body movements?
- Q4. How easy was it to move while wearing the belt?
- Q5. How unobtrusive was the belt?
- Q6. How comfortable was the belt?
- Q7. How ergonomic was the belt?
- Q8. How lightweight was the belt?
- Q9. How well did the belt fit your waist size?
- Q10. How would you rate the belt's physical discreetness?
- Q11. How silent were the belt vibration motors?

The functionality and performance questions for the researcher were:

- Q1. How easy was it to create your desired configuration of the belt, which involved adding/remove factors, moving factors around, wiring, etc. (take into account scalability and reconfigurability)?
- Q2. How easy was it to design your desired vibrotactile patterns using the GUI (take into account the expressiveness of the system, and GUI usability)?
- Q3. Was the portability of the belt, in terms of wearability and wireless capabilities, suitable for your intended application?
- Q4. Was the durability of the belt suitable for your intended application?
- Q5. Was the reliability of the belt suitable for your intended application?

Q6. Was the wireless communication latency suitable for your intended application?

Q7. Was the battery life suitable for your intended application?

In order to determine how well the experiment itself fared, we devised 5 research hypotheses for objective evaluation:

Q1. Subjects will achieve at least 90

Q2. Subjects will achieve at least 85

Q3. Subjects will achieve at least 85

Q4. No one spatio-temporal pattern is more difficult to identify than the other.

Q5. The moves of neither dance-modified box step or modified electric slide-will be more difficult to recognize than the other.

Other than the objective evaluations, participants were asked questions directed towards the dance experiment itself:

Q1. How easy was it to recognize the vibrotactile patterns?

Q2. How intuitive was the mapping between vibrotactile patterns and movements you had to perform?

Q3. In the second part of the testing phase, you learned how to perform a dance sequence. How well did you learn the dance through use of the vibrotactile patterns?

Q4. If you wanted to learn how to dance someday, how likely are you to use this system?

Q5. Do you think others would like to use this system to learn dance?

Q6. Have you danced before?

Q7. If you have danced before, how many years?

Q8. What is your preferred style of dance?

10.10.4.5 Results

Usability: In order to understand the usability of the haptic belt through the subjective evaluation survey, we performed a one-way ANOVA on the data presented in Fig. 10.16. Considering a 5% significance test on the null-hypothesis that there is no significant difference in the means of the 11 usability questions, a 10 DoF along the questions axis, and $11*(13-1) = 132$ DoF along the participant axis, F test results in $[F(10,132)=3.29, p=0.0008]$, thereby rejecting the null hypothesis. Further, as a post-hoc analysis, a Multiple Comparison Procedure on the linear one-way ANOVA (with significance level $\alpha=0.05$) shows that with respect to question 2 (How easy was it to take off the belt?) and question 4 (How easy was it to move while wearing the belt?), question 10 (How would you rate the belt's physical discreteness?) and question 11 (How silent were the belt vibration motors?) are significantly different, thereby contributing to the rejection of the null hypothesis. On reviewing the descriptive evaluation provided by the participants on the haptic belt, it was discovered that question 10 relating to the physical discreteness was rated low due to the bulkiness of controller box on the belt. For question 11, referring to the noise made by the vibration motors, a redesign of the tactor modules is necessary to ensure that the vibration motors are enclosed rigidly within the tactor module.

Question 4, relating to how easy it was to move wearing the belt, had the highest mean value of 9.46 (SD 1.13). This question relates to the important aspect of whether the belt allows the participants to move freely wearing the device. Any research platform has to offer this movement flexibility so that the hindrance due to the platform does not bias the user's opinion of the experiment's research questions. It was also seen that it was easy to take off the belt (Question 2) when compared to putting it on and adjusting the location of the vibrators (Question 1). The results are obvious as removing the belt necessitated only releasing the plastic snap buckle, whereas, wearing the belt and locating the motors necessitated the participant's attention and effort.

Functionality and Performance: Fig. 10.17 shows the responses of the independent

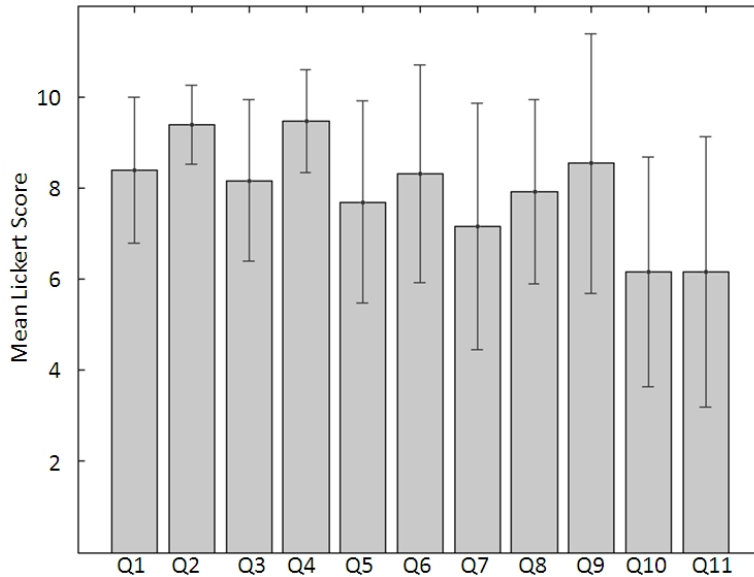


Figure 10.16: Usability Results.

researcher to the seven questions on functionality and performance. Since the belt was reviewed by one independent researcher, no formal statistical analysis can be done on the results. We report here our observations from what the researcher offered as explanations to his survey. No problems were experienced by the researcher when reconfiguring the belt or designing spatio-temporal patterns. In terms of the performance of the belt, portability, durability and wireless communication latency were found to be fine. As can be seen from Fig. 10, the two important drawbacks in terms of functionality and performance were found in a) the reliability of the belt for the intended application (Question 5) and b) the battery life of the haptic belt (Question 7). The failure to meet the necessary battery life on the belt was realized by the developers through the experimental study itself. The choice of battery manufacturer turned out to be a problem and has little or nothing to do with the design of the power supply module for the belt. We also realized that the researcher found the battery issue to be the main reason to consider the reliability of the belt to be low or not up to expectation.

Quantitative Evaluation of the Dance Experiment: Fig. 10.18 shows participants' recognition accuracies on the 12 spatio-temporal patterns that were delivered as part of

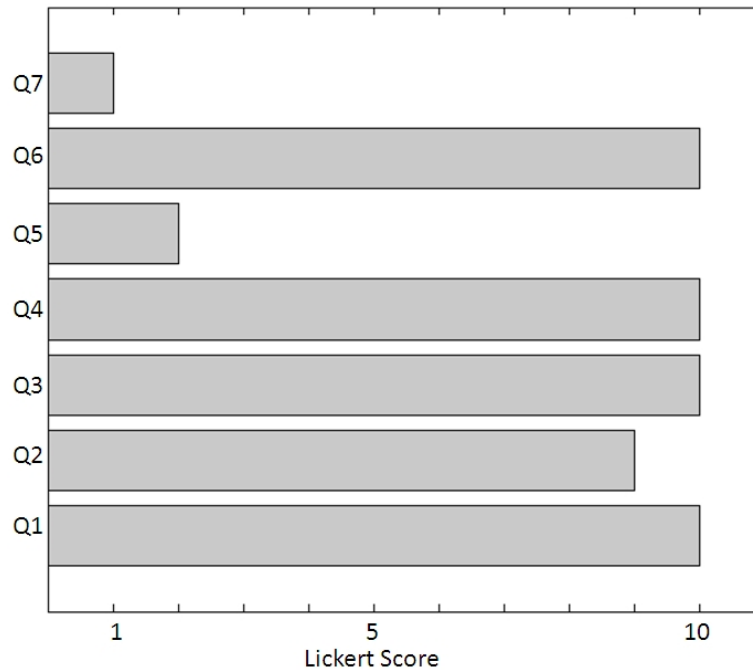


Figure 10.17: Functionality and Performance Results.

the dance experiment. The overall recognition accuracy of vibrotactile patterns, averaged across participants, was 97% (SD: 4.6%). The average accuracy for recognizing the individual moves of the modified box step dance was 88% (SD: 20%), and the average accuracy for recognizing the individual moves of the modified electric slide was 95% (SD: 7.5%). Fig. 11 shows the results of the experiment where the participants performed the 12 patterns based on the 12 spatio temporal sequences. These results support hypothesis (1), showing that overall, participants had no difficulty recognizing the vibrotactile patterns. Using a one-way ANOVA, no significant difference [$F=1.87$, $p=0.0475$] between average recognition accuracies of vibrotactile patterns was found. This supports hypothesis (4), and shows that no single pattern was more difficult to recognize compared to the others. These results also support hypotheses (2) and (3), showing that participants were able to link moves together to perform some basic dances. A one-way ANOVA was applied to the accuracies achieved on the two dances, revealing no significant difference [$F=1.55$, $p=0.2255$] between the average recognition accuracies of the two dances (modified box step and modified electric slide). This supports hypothesis (5), and shows that participants didn't find one dance

more difficult than the other, even though the electric slide is longer and more complex than the box step. However, 3 out of the 13 participants scored very low on the modified box step dance, after which they performed very well on the more complex electric slide dance. We hypothesize that, for these participants, additional learning beyond the familiarization phase was required to learn how to link movements together; we believe that this learning took place during the modified box step dance steps. Reversing the dance sequences may have avoided this, but we feel that performing box step dance before the electric slide dance facilitated learning, as the box step dance is easier than the electric slide.

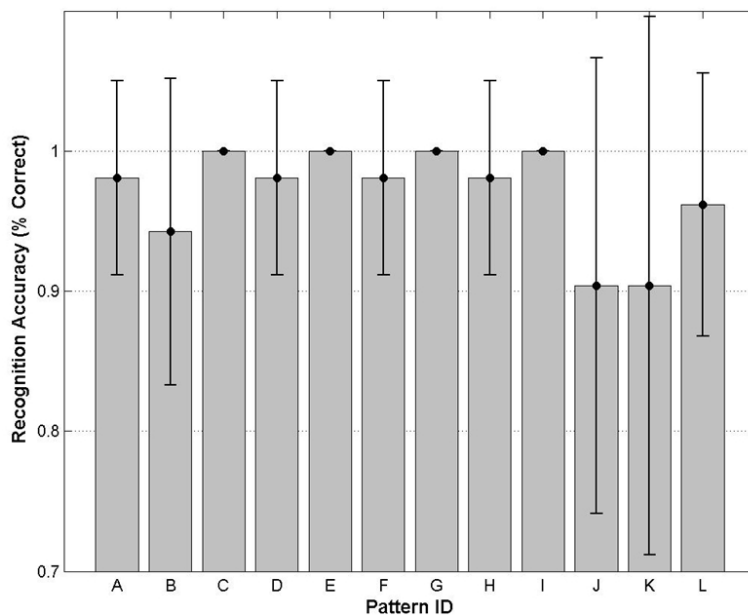


Figure 10.18: Pattern Recognition Results for Dance Experiment.

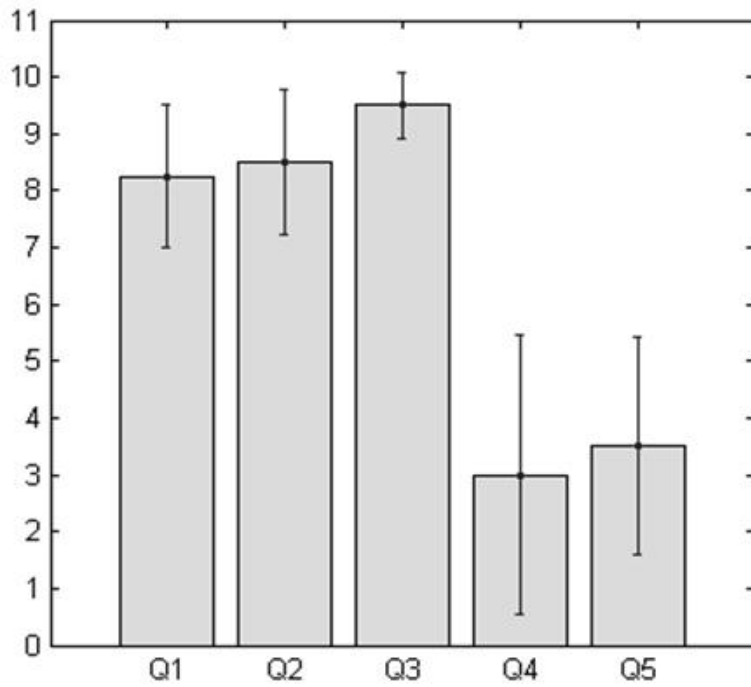
Subjective Evaluation of the Dance Experiment: Fig. 10.19 shows the subjective user responses for questions 1 through 5 based on whether the participants were experienced in dancing or not. Questions 6 through 8 explored the participants dance experience level, and we found on average, participants had no experience with dancing to about 5 years. We set the average of all user experience (1.8 years) as a threshold to decide whether participants were experienced or not.

Figure 10.19(a) shows the results of participants who were experienced (mean experience of 5.12 years), and (b) shows the results of participants who were inexperienced

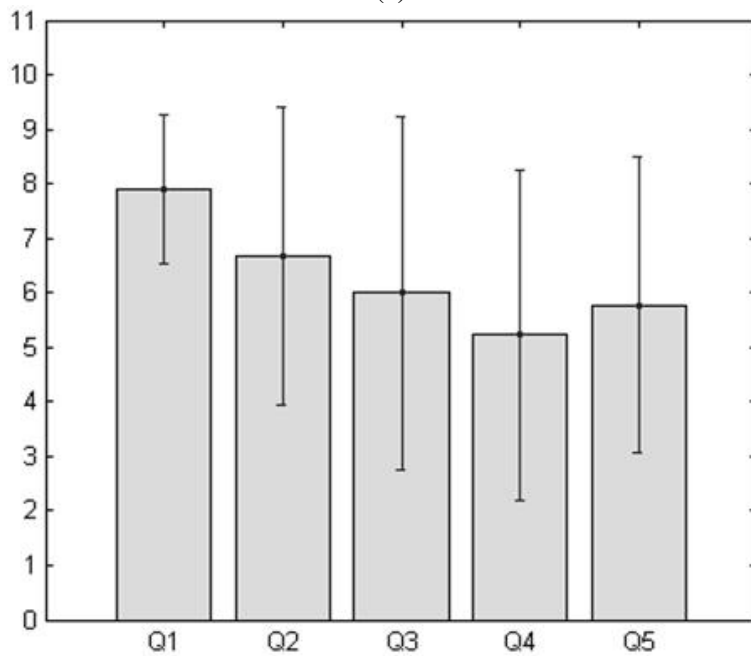
(mean experience 0.44 years). From Figure 10.19, it can be seen that the participants' opinions varies widely between the experienced and inexperienced groups, except for question 1, which enquired about the ease of recognizing the spatio-temporal patterns. The mean response for question 1 was 8 (SD 1.16). When the participants were asked how intuitive (Question 2) and useful (Question 3) the spatio-temporal patterns were, the experienced group seemed to desire having this sensory augmentation more than the inexperienced group. Correlating this to the quantitative analysis, the experienced group achieved 99% accuracy (SD: 1.8%) in recognizing all the 12 spatio-temporal patterns, whereas, the inexperienced participants achieved 95% accuracy (SD: 5.2%). We hypothesize that the experienced dancers had no problem executing the dance step and hence could focus on the vibrotactile pattern, whereas the inexperienced participants had to consciously process the haptic cues and the movements. When the participants were asked how likely they would use this device again (Question 4), or suggest this device to someone else (Question 5) to learn dance, the results seem to indicate opposite of what was seen in the previous two questions. The experienced dancers found this device rudimentary and not recommendable, whereas the inexperienced dancers seem to reluctantly agree to using or suggesting a sensory augmentation.

10.11 Group Interaction Dynamics Through Haptic Belt

As mentioned above, the haptic belt forms a situational awareness system that is capable of delivering not only the positional information of the interaction partners, but also their dynamic movement in the environment in front of the user. This in turn would help the people who are blind and visually impaired to understand the social scene in front of them. In the following section, we show the social interaction assistant in its entirety for assisting people who are blind and visually impaired in group interactions and dyadic interactions.



(a)



(b)

Figure 10.19: Questionnaire Results from Dance Experiment for Experienced Dance Participants (a) and Inexperienced Dance Participants (b). Responses from Q6-Q8 are excluded.

THE SOCIAL INTERACTION ASSISTANT & ITS ROLE IN SOCIETY

This dissertation introduces the concept of social situational awareness and discussed in detail how to develop a social interaction assistant for aiding people who are blind and visually impaired in everyday social interactions, through the use of a evidence-based modeling of social needs of the target user population. Accordingly, we identified two important areas of social interaction assistance, namely, the dyadic interaction assistance and the group interaction assistance, while also providing a technology to provide immediate intervention into any asocial stereotypic body mannerisms that they individuals might be displaying.

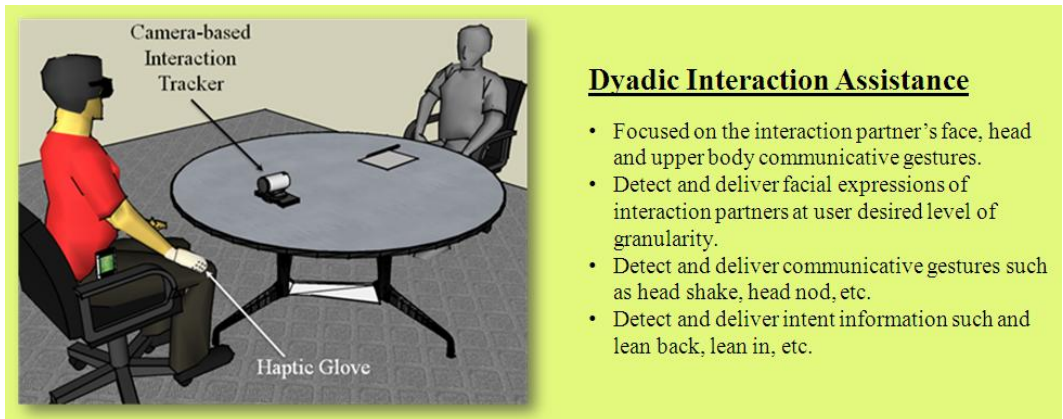
Figure 11.1 shows the Dyadic Interaction Assistant and the key highlights of the proposed solution. Figure 11.1(a) highlights the use case scenario where the device consists of a camera that is placed in front of the user and it pointed in the direction of the interaction partner. The user in turn uses the Haptic Glove to receive all forms of interpretations of the interaction partner's face on the back of the hands. The camera used in this application uses a micro servo mechanism to keep the face of the interaction partner in focus and always centered within the video stream generated. This helps in ensuring that the facial expressions are never missed in during the interactions. As part of the future work, the user will be provided with a mechanism for controlling the amount of information received from the interaction partner. The interface used by the user is highlighted as part of the Group Interaction Assistant in Figure 11.2(b), termed as the Clicker. The control interface uses a series of buttons using which the user can choose four levels of details of the interaction partner's face, namely,

Literal The facial mannerisms are literally translated to the fingers as vibrations, with the system making no decisions on the nature of the facial movements.

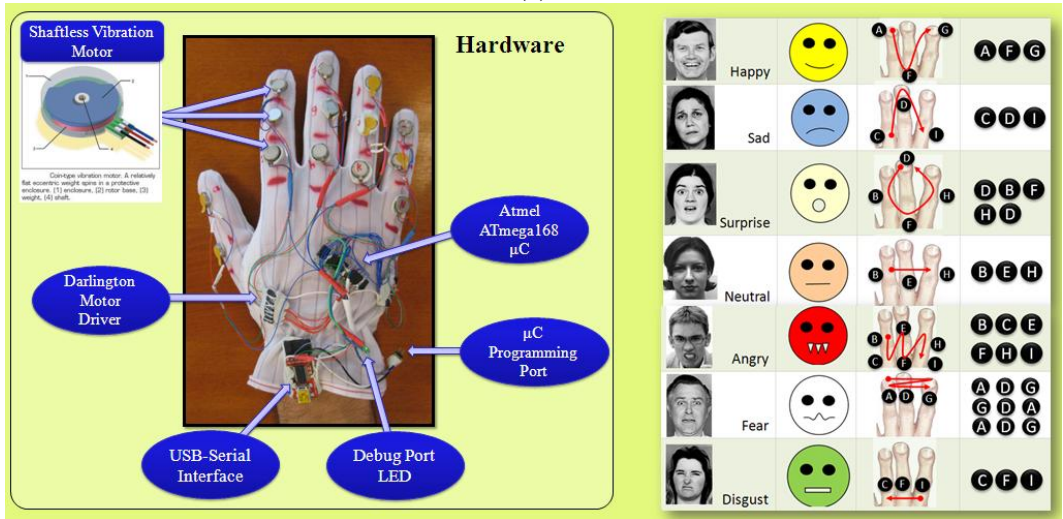
Symbolic The facial movements are classified into groups that are previously chosen, such as happy, sad, surprise, etc. (Example of 7 classes shown in Figure 11.1(b)).

Semi-Literal One step above the raw transfer of facial mannerisms. This could include modeling the mouth as curves instead of simple points.

Semi-Symbolic One step below Symbolic representation, this level would provide higher level information than just expression decisions, like the Action Units on the face.



(a)



(b)

Figure 11.1: Group Interaction Assistance; (a) Scenario for group interaction assistance, (b) The integrated group interaction assistant.

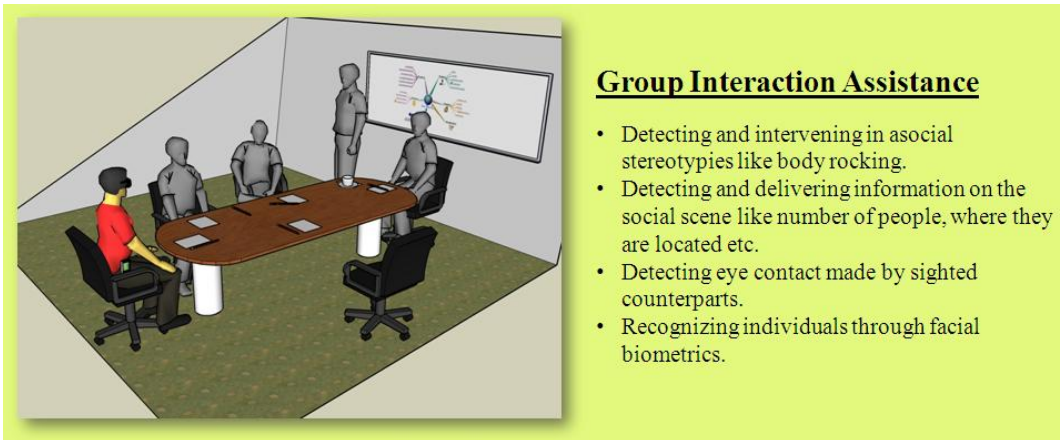
Figure 11.2 shows the Group Interaction Assistant as described in the chapters above. Figure 11.2(a) shows a use case scenario where the user who is blind is facing a group of people in a meeting. The group interaction assistant gives some of the important cues such as the position of the people, their movement around the user, eye gaze of the

individuals, etc. The primary sensing and deliver mechanisms of the group interaction assistant are shown the Figure 11.2(b). The sensing is based on the glasses that is worn by the user with a camera mounted on the nose bridge of the glasses. The delivery is primarily centered around the waist belt that acts as the haptic actuator. Along with the sensing and delivering the proxemics information, the group interaction assistant also provides the user with information about his or her own body mannerisms. To this end, the device uses a motion sensor that is placed just below the neck of the user. This sensor is capable of picking up body movements and in this application, focused on detecting body rocking. The feedback could come in terms of the vibrations or using a earphone placed behind the ear lobe so it does not affect the user's hearing.

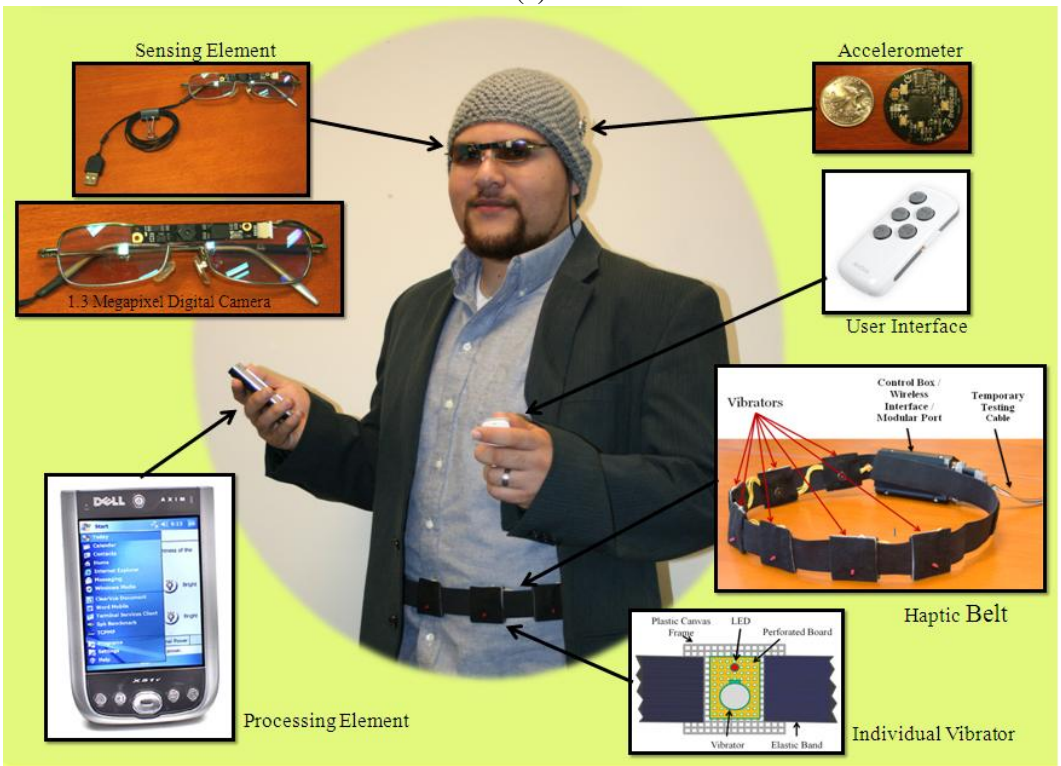
An important component of the group interaction assistant is the integration of the facial biometrics into the proposed technology. This allows the user to not only know where people are in front of them, but also know who exactly is located where. In future, we propose to integrate the group and dyadic interaction assistants, there by allowing the user to know exactly where any specific person is located and be able to focus the camera towards that individual to receive their facial expression cues. To this end, we are developing a back worn haptic face display that will be able to integrate into a jacket. This will allow the user to feel the vibrations on their back without having to wear a glove explicitly. The controller (Clicker) acts are the primary interface for the user to enable or disable features on the device. Once the integration of the dyadic and group interaction assistants is complete, the user should be able to surf their field of view and assess each individual intricately.

11.1 Role of Social Assistance in the Society: A Policy Discussion

Having introduced a novel concept in assistive technologies, namely the social interaction assistant, we would like to spend sometime discussing the role of such a technology in the society. The study of interactions of technology and society provides an opportunity to understand the real contribution of a certain technology or concept to the economic, social and cultural future of a society. Especially since the proposed social interaction assistant bears its roots in the social interaction issues that form the core of human societies, we find it in-



(a)



(b)

Figure 11.2: Group Interaction Assistance; (a) Scenario for group interaction assistance, (b) The integrated group interaction assistant.

triguing to understand not only the broader impact of such a technology, but the need for it in general. Further, recently, there has been a growing awareness on the role of technology on society. This growth in interest can be correlated to various technologies that have unforeseen consequences on the public, which could have been averted if the artifacts of new technology could have been studied before their mass proliferation. It is always impossible to demonstrate the role a certain technology will assume as the parameters involved in assessing the consequences are innumerable, the most problematic of which are the humans and their behavior towards a certain technology. A clear example of such a technology in the recent decades is the prolific use of cameras in every form of technology from surveillance systems to reverse view cameras in cars. While the clear advantage of the technology is demonstrated by its use, there are unprecedented privacy issues associated with the use of cameras. Such privacy issues are heightened by the proposed technology where we consider the use of wearable cameras on the user who is blind or visually impaired. While it can be argued that the use of camera is restricted to the user alone and that no explicit user images will not be saved on the device, it is impossible to inform all interacting partners that there is a camera on the user and that they are being captured. Before getting into the technology debates related to wearable cameras, we highlight the need for social assistance in the society and make a policy case for the proposed social interaction assistant. Following the social impact arguments, we discuss the implications of wearable cameras in public space.

11.1.1 Social Disability: The hidden barrier to professional growth in the disabled population

On February 12, 2009, Vice President Joe Biden announced the appointment of a Special Assistant to the President for Disability Policy. By selecting the Associate Director of White House Office of Public Engagement, Kareem Dale, to this post Obama became the first President in the US history to have a special policy advisor overseeing disability issues.

One of the key components of Obama's policy on disability is the effort to increase access to employment for the disabled population. The White House website on Disability

¹ states:

President Obama is committed to expanding access to employment by having the federal government lead by example in hiring people with disabilities; enforcing existing laws; providing technical assistance and information on accommodations for people with disabilities; removing barriers to work; and identifying and removing barriers to employment that people with public benefits encounter.

With this statement the Obama administration has identified some important employment issues associated with disability, but it seems to be focused only on bringing people with disabilities into the labor force. Once they are in, people with disabilities face a number of additional hurdles. The barriers identified by the administration fall short of addressing how to retain this population in the work force and enable a self-propelled professional growth.

What does it take to succeed in one's career? Evidence shows that the social skills of individuals and their ability to integrate themselves into the work environment are incredibly important. Social skills - such as making friends at workplace, ability to lead a team, facilitating decision making in large teams, conveying confidence, etc. - all play a vital role in sustained professional growth. Unfortunately, people who are severely disabled, like those who are blind, often find it difficult to assimilate into the social atmosphere of their work place with the same ease of their functionally able counterpart.

“There is no professional growth without social skills”, explains Dr. Terri Hedgpeth, director of Disabilities Resource Center on Arizona State University's campus. Hedgpeth has been blind her whole life and had to learn how to socially interact with her sighted colleagues. For instance, she learned to turn her head towards her interaction partner to mimic eye contact. She learned to hear people's bodily movements to assess what they were communicating non-verbally.

Hedgpeth doesn't want to be offered any social leeway because of her disability,

¹<http://www.whitehouse.gov/issues/disabilities/> extracted Dec 11, 2010

but strongly believes in training individuals who are blind and visually impaired to learn the same social skills as their sighted peers. She laments the fact that currently there are no federal programs, either vocational or academic, that trains people who are visually impaired about important social skills in professional setting. Social training is generally reserved for children and young adults who have a severe case of tics, like body rocking or eye poking.

11.1.2 Effect on Social & Emotional Intelligence due to Disability

The social disconnect is not limited to visual impairment and blindness alone. Social implications of disabilities can be seen across the spectrum from physical disabilities like wheelchair and quadriplegia to cognitive disabilities like Autism. Studies in Cognitive Psychology support the hypothesis that social interactions play a vital role in the overall development of overall intelligence in humans, especially, in the development of Social Intelligence (or Interpersonal Intelligence [374]) and Emotional Intelligence [375]. Social and Emotional intelligence are vital components in an individual understanding the importance of other people and things in their surroundings. Without active social interactions, a large part of the learning component is lost.

First defined by Edward Thorndike, Social Intelligence is "the ability to understand and manage men and women, boys and girls, to act wisely in human relations" [376]. Karl Albrecht [377] argues that Social Intelligence is the basis for five important aspects for an individual to mingle into his/her society, including, 1) Situational awareness, 2) Sense of Presence, 3) Authenticity (or Individuality), 4) Clarity (of action), and 5) Empathy.

Along similar lines as social intelligence, neuro-psychologists have defined Emotional Intelligence (EI) as the ability, capacity, and skill to identify, assess and manage the emotions of one's self, others and of groups of individuals. If disabilities, like visual disability, becomes a barrier to seeing the emotion artifacts, the individuals might find themselves lacking the ability to interact with their society at a level that is otherwise considered normal. Many models have been proposed in the past to explain EI, such as Ability based

models [378], Mixed models [379] and Trait based models [380] and all these models point towards the fact that reduction in social interactions can reduce the overall understanding of an individual of their place in the society. Recently, EI metric scales have been used to diagnose autism spectrum disorders, including autism and Asperger syndrome, semantic pragmatic disorder or SPD, schizophrenia, and Attention-deficit hyperactivity disorder (ADHD). These measurements have shown a direct correlation of one's ability to increase their overall emotional involvement within the society by increasing their social interactions.

Primate researcher, Humphrey [381], has demonstrated the real-world effect of social interactions to cultural transmission of knowledge and the development of intelligence. His studies with rhesus monkey have emphasized the positive influence of social interactions on the development of general intelligence. For example, Helen (a rhesus monkey) had her visual cortex surgically removed and studies were conducted on her recovery of spatial vision. Over four years, isolated within the laboratory, Helen hardly recovered any of her spatial knowledge. However, when she was taken out of the laboratory into the real world and allowed to interact with objects and other monkeys, she regained three dimensional spatial vision within a few weeks. Humphrey argues that the interactions with other monkeys were key to Helen's learning of interactions (both with objects and other monkeys). Could this be true with humans also?

From a neuro-physiological perspective, advanced functional brain imaging is enabling researchers to study the workings of human brain under various functional conditions and they are confirming the role of social intelligence as an important aspect of human learning. Brothers [382] has worked extensively on the neuro-physiological patterns in primate brains that are associated with social behavior. Her work has established the presence of dedicated brain regions involved only in social cognition (Social cognition is the processing of information that culminates in the accurate perception of dispositions and intentions of other people). She has proposed a network of neural regions that comprise the social brain and she argues that a malfunction of the any component of the social brain results in

reduced social cognition. Her work has been recently bolstered by [19], where the authors study autistics and controls under functional Magnetic Resonance Imaging (fMRI). The subjects watched another person's eye expressions, and guessed what that person was thinking or feeling. The fMRI images confirmed Brothers observations of STG and amygdala activations during social cognition, and showed that people with autism display a cognitive disability in the amygdala which prevents them from making appropriate mental inferences of other people's emotions or facial expressions. Authors conclude that a social brain does exist, and that teaching children and adults social skills could offer a means of increasing activations in the social brain. This conclusion is supported by the behavioral research in autism that employs social interaction training and language skill training in children to ameliorate the social deficits characteristic of autism spectrum disorders (ASD).

The disabled population faces social barriers due to their sensory, motor or cognitive dysfunction. Overcoming this social barrier cannot happen overnight through "enforcement of laws" as reported on the White House's disability policy website. This has to happen through strong dedicated programs that study the social barrier to employment in the disabled population and offer effective solutions (social assistive aid, social education programs and co-ed of disabled and non-disabled children to encourage mutual learning of social skills, to name a few) to reduce the consequences of social disability.

11.1.3 Psychological Breakdown related to Social Skills

Recent studies by Segrin et al. [72] have shown that poor social skills are antecedents to psychosocial problems including depression, loneliness, social anxiety, etc. The authors conducted a battery of tests on college students to determine the effect of stress on the students when they live at away from home. Figure 1 shows Depression and Loneliness plotted against stress levels of undergraduate students. Depression was measured using the Beck Depression Inventory [383], while Loneliness was measured on the UCLA Loneliness Scale (version 3) [384] as an index into the students experience of loneliness. For both of these tests, the participating students were categorized into high, medium and low social skilled groups based on the Social Skills Inventory [385] (a battery of tests administered to

determine the socialization ability of an individual).

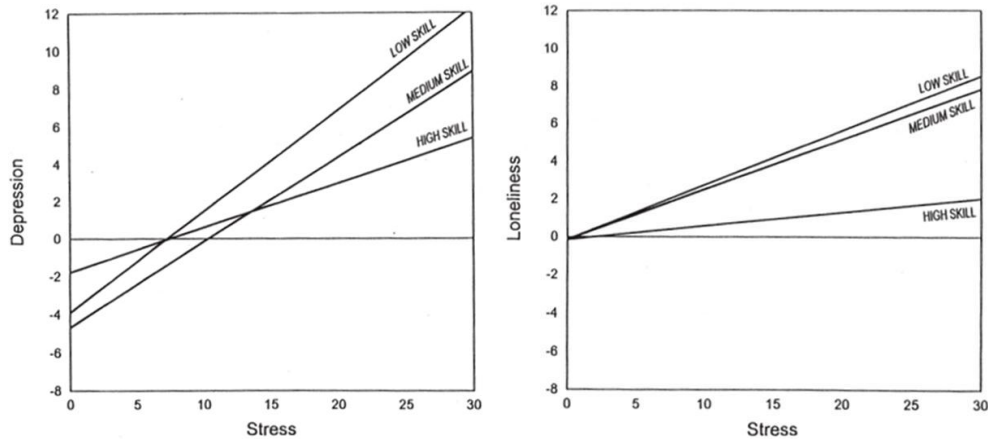


Figure 11.3: Depression and Loneliness of students plotted against stress levels in high, medium and low social skilled undergraduate students. (Please see text for the scales used for the measurement.)

One can immediately identify a positive correlation between stress and an increased experience of psychosocial problems in all the students, but the ones that rank higher on social skills show higher resistance to stress and in turn higher resistance to mental breakdown. Students assessed with mild or lesser social skills were highly vulnerable to social issues as the stress increased.

Similar results were found in [3] where the authors conclude that people with high competence in communication are known to display immense capability towards adapting their social behavior based on others in their surroundings. Such competence has been acknowledged to reinforce social skills thereby creating a reinforcement feedback that allows these individuals to be successful in their social endeavors [386] and in turn successful in their life. In a tangential study, though Magnusson [387] was not looking for social interaction needs in people, found that social interaction is an important dimension in the cognitive organization of human behaviors. When college students were assessed individually, and as a group, to determine how they classified everyday activities into different situations, Magnusson discovered 5 dimensions (Principle Dimensions). These included two dimensions based on whether the students perceived a situation as being positive (positivity) or negative (negativity) influence on their behavior, two dimensions based on whether the situ-

ations were active or passive, and finally, the fifth dimension was based on social interaction with others. His study emphasizes how social interactions are perceived by individuals as an important scale for judgments on their activity of daily living.

11.1.4 Societal Metrics of Social Disability

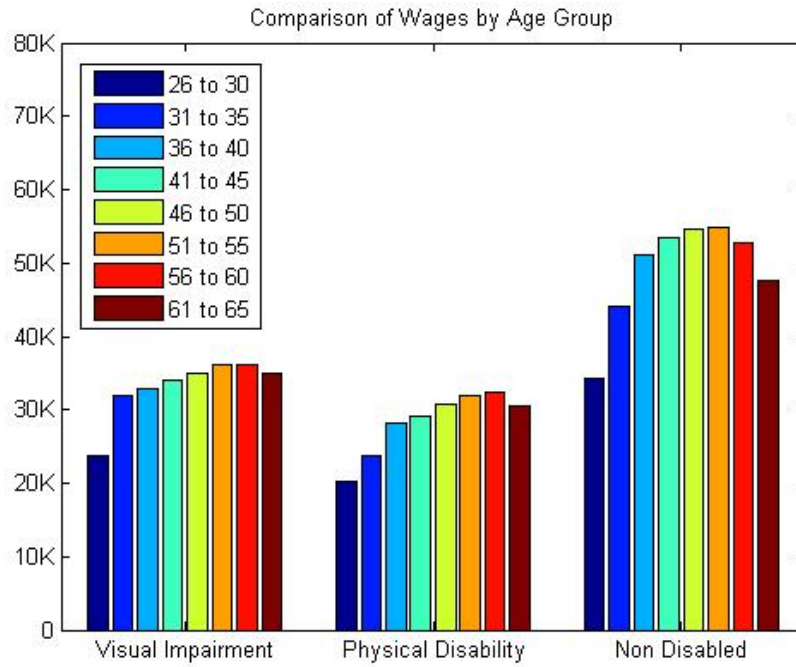
11.1.4.1 Comparison of Economic Status

The Americans with Disabilities Act was signed into law in 1990 by President George H. W. Bush and from then on, pretty much the same verbiage has followed each successive president's agenda on disability employment. As a nation, we have been successful in moving this segment of population into the work force, but still there is a large gap in their professional success. On a wage comparison scale (data extracted from the 2008 American Community Survey (ACS) questionnaires²), the US visually impaired population make on an average 32% less than general population of the same age. The US physically challenged population makes 42% less. See Figure 11.4. When one includes education level in these statistics, the results are even more disappointing. People with visual disability, with post graduate education, make 47% less than the average population with post graduate education and people with physical disability with post graduate degrees make 72% less than the general population with similar degrees.

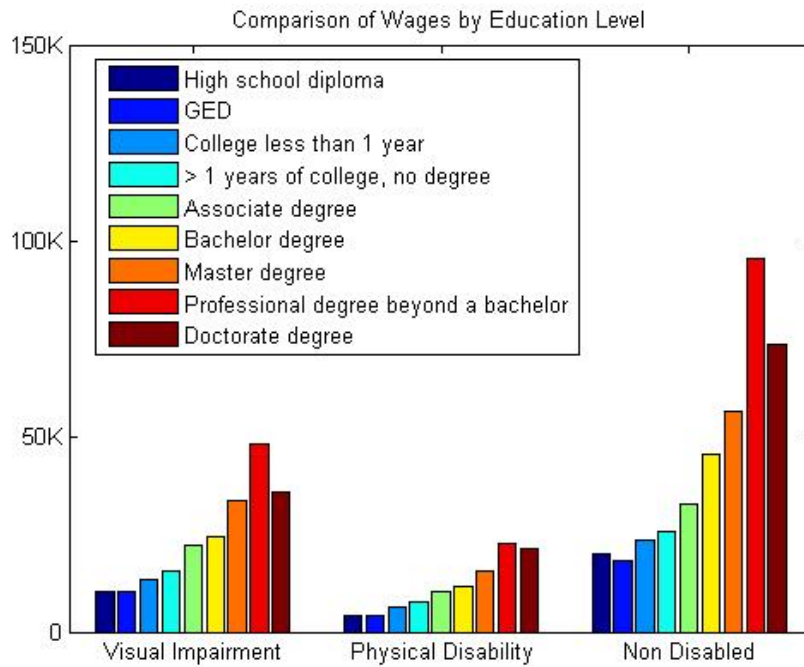
11.1.4.2 Comparison of Personal Lives

Similar to the wage comparison, we also focused on the personal lives of individuals with disability and compared it with the personal lives of individuals without any disability. The goal was to compare the companionship among the various populations and decipher any correlations that may exist between disability and the number of times people got into and out of civil unions. Hedgpeth explains that disability could come in the way of so-

²The ACS is a long format survey conducted on a yearly basis with the American public. The survey is carried out in a format similar to the American decennial census survey. Current survey standard seeks 10% of the population (3 million individuals) and are requested to respond to a long format survey that asks questions on the living standards, salary, family structure, housing quality etc. A long format of the survey can be found in the Appendix C. The statistics used for the analysis here has 27990 individuals who were visually impaired alone with no other associated disabilities, 155726 individuals with only physical disability and no other disabilities, and 2591521 with no disabilities whatsoever.

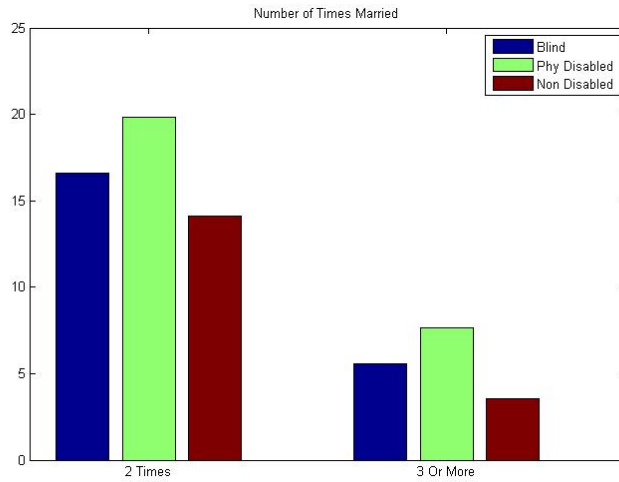


(a)

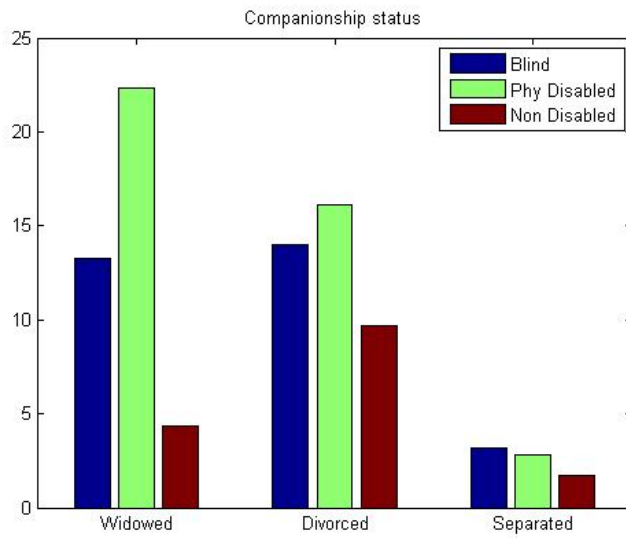


(b)

Figure 11.4: Comparison of annual wage of the visually impaired, physically disabled and non-disabled population. (a) Compared by age group. (b) Compared by education level.



(a)



(b)

Figure 11.5: Comparison of personal lives of visually impaired, physically disabled and non-disabled population. (a) Number of times married. (b) Companionship status.

cial personal interactions and many times disabilities could impose unforeseen pressure on the marriage. Figure 11.5(a) shows a comparison of the number of times individuals were married, as compared between visually disabled, physically disabled and the non-disabled population. On an average the number of marriages is higher among the disabled population and specifically higher among the physically disabled. Similar results are seen while comparing the companionship status of individuals who are disabled. There is a higher chance of widowhood, divorce and separation among the disabled populations when compared with the non-disabled population.

Hedgpeth emphasizes the point that people skills are the most important tool for professional success. Unfortunately, social disconnect is a repercussion of disabilities. It is important to train the disabled population to circumvent their social disconnect, while we train the rest of the population to understand and acknowledge this social barrier. To this end, some of the concepts introduced as part of the social interaction assistant proposes to address them by providing access to the social cues that are otherwise inaccessible to the individuals who are disabled. While the discussions in this dissertation are restricted to sensory disability, specifically vision disability, and discuss all the social problems associated with it, the same could be said of the problems faced by people with other sensory disability (like hearing loss). We reserve our discussions only to sensory disabilities as the problem of determining appropriate adaptations for overcoming cognitive disabilities require research that is beyond the scope of this document.

11.2 Wearable Cameras: Ethics & Privacy

Wearable cameras and cameras on smart phones have started a debate on their ethical use and the possible misuse through infringement of privacy. As discussed earlier, the social interaction assistant discussed in this dissertation poses a significant challenge in terms of the ethical use and privacy as the camera used on the device is always recording (a term that needs to be considered carefully based on the underlying technology) the interactions of the user and their interaction partner. The discussion of wearable cameras and their role in the society has been debated by a number of related areas including, privacy proponents,

civil liberties union, law makers, law enforcement, security and surveillance experts and finally the users. A very good example of a comprehensive discussion on the role of wearable cameras on the society can be found in [388]. Unfortunately, most of the times these discussions happen after a certain technology takes root in the society and most of the times it seems to be a retroactive fix to the privacy and ethics issue. In contrast, we would like to bring these topics to the forefront of the technology development to ensure that the device does not violate the ethics or privacy issues as far as possible. While the role a certain technology will play in the society can never be predicted with certainty before it is introduced and certain level of acceptance has come into play, it is possible to direct the development and deployment of the technology in directions that could result in fewer ethical barriers.

Morgan and Newton, in their discussions on the issue of *Protecting Public Anonymity*, highlight an interesting exercise carried out at public policy class in Carnegie Mellon University. We reproduce their writings from [389] here for clarity on their exercise.

... For years, we have used a teaching case in Carnegie Mellons Department of Engineering and Public Policy, in which graduate students are asked to assume that a basic smart car system is about to be implemented. They are asked to consider whether the state should run a pilot study that would implement a number of advanced system functions, such as insurance rates that are based on actual driving patterns, externality taxes for congestion and air pollution, and a system for vehicle location in the event of accident or theft. We find that students immediately assume a system architecture that includes real-time telemetering of all vehicle data to some central data repository. Then they become deeply concerned about issues of civil liberty, invasion of privacy, and social control, and often go on to construct arguments that such applications should be banned.

It is often not until students have worked on the problem for several hours that someone finally stumbles on the insight that most of the difficulties they are concerned about result from the default assumptions they have made about the systems architecture. If information about vehicle location and driving performance is not telemetered off vehicles on a real-time basis, but is instead kept on the vehicle, not as a time series but in the form of a set of simple integral measures (such as a histogram of speeds driven over the past six months), then insurance companies could access it twice a year with all the time resolution they need. If detailed records of who drove where and when are not created, then most of the civil liberty problems are eliminated. Many of the potential concerns raised by other system functions in this teaching case can also be largely or entirely eliminated through careful system design choices.

This example very well exemplifies the problems that technologists get into, if proper care is not taken towards designing the system from ground up. Similar to the case above, the social interaction assistant would really benefit from collecting all of the interaction data from the users and using them for further machine learning and signal processing. Unfortunately, this is not a viable solution as recording of video, especially of others, without their knowledge could result in privacy violations. To this end, we highlight on various steps that could be taken towards ensuring that the technology does not *record* person identifiable data that have immediate consequence to a person or institution who are not interested the application of social assistance. For example, consider face recognition as a biometric, institutions are fighting the problem of securely storing the face images of their employees or members who are allowed into a building. The loss of the face images could immediately affect privacy of these individuals because images are human readable and have applications in various areas. On the other hand, if instead of saving the human readable face image, a certain proprietary encoding was applied to the face image and only the resulting features (a machine learning terminology applied to the end result of a transformed data) are saved, any loss of data could then be cordoned off from abuse by securing or destroying the methods used to create the features.

We highlight here some concepts that could potentially reduce the problem of ethical and privacy concerns associated with the social interaction assistant.

- *Never saving human readable data* As explained in the earlier section, with the current state-of-the-art in data representation technologies, it is possible to transform the personal data from the users into features that only proprietary algorithms could decode. By isolating the algorithm from the collected data, it should be possible to secure the data from ever being available in a human readable format.
- *Erasable Biometrics* Erasable Biometrics is a newly emerging area of biometrics where it is essential, by law, to remove all biometric data collected on an individual when once he/she is not part of an institution that collected any such data. Pioneering work in the area of biometrics use a data morphing model as shown in the Figure 11.6. When an individual is added to an existing biometric database, he/she is assigned a private key which is user-specific and can be destroyed if need be. All the biometric data collected from the user will then be transformed from a global set of transformation parameters, which are then randomly encoded based on the private key. All machine learning takes place on top of the transformed biometric data and never on the actual data which is discarded immediately after its first use for registering the user. When once a user prefers his data be removed from their system, the erasing happens by destroying the user supplied private key. Since the key controls all the transformation parameters, there is no way the transformed biometric data be ever brought back into the human readable insecure format.

Having introduced Erasable Biometrics, it is possible that such techniques could be adopted to ensure that any individual who prefers at a later time to be removed from the social interaction assistant could then be considered erasable.

- *Design intelligent sharing of data* Uncontrolled sharing of digital data is becoming a menace to the future of our society. Users prefer a certain level of sharing that is impossible to prevent and at the same time care has to be taken to ensure that important private information is not lost. To this end it is essential to design intelligent

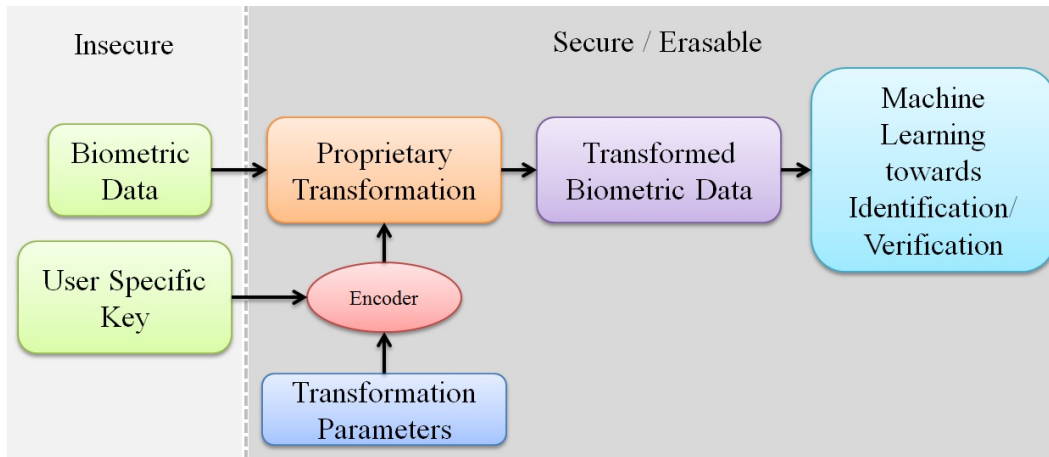


Figure 11.6: Framework for Erasable Biometrics

user sharing where information on interaction partners could be shared among users who are blind.

- *Other additional security features* In an emerging society of digital data sources, security has emerged as an important issue and everyday security features like password blocking, auto lock of device activity, theft control through tracking, etc. should be considered as an essential component of the social interaction assistant design along with the other constraints that were identified in Chapter 3. See [389] for an exhaustive list of specific design criteria for preventing collection of user-specific data from devices.

In conclusion, while social interaction assistant is well placed to make a tremendous impact on the quality of life for individuals who are blind and visually impaired, there could be possible ethical and privacy issues associated with the wearable camera on the user. (Not to mention the awkwardness a user may have to face if he/she had to wear a t-shirt all the time informing the interaction partner that they are being recorded.) This chapter tried to highlight some of the important privacy issues and proposed solutions that could potentially overcome some of them. As discussed earlier, it is very difficult to determine all the impact that a certain technology will have on the society before it is introduced into the user space,

but it is always possible to hypothesize the potential risks that the designers may expose to the society to. We have attempted such an exercise in the above discussion.

Chapter 12

Conclusions & Future Work

In this dissertation, an evidence-based methodology towards enhancing social mediation between individuals is presented. Specifically, the social mediation technologies described here attempts to provide people who are blind and visually impaired to receive non-verbal communication cues from their sighted counterparts. Chapter 2 discussed the need for enriching social situational awareness in everyday personal and professional lives of individuals. Chapter 3 highlighted the importance of enriching social situational awareness for individuals who are blind and visually impaired and lays foundation for the bulk of the work presented in this dissertation. Chapter 4, 5, 7 and 8 discussed various technologies that can enable users who are blind and visually impaired to access social signals that are important for having a rewarding social interaction with sighted counterparts. The details of these chapters will become clear once the reader has been introduced to the various social situations that require attention, as detailed in Chapter 3. Chapter 6 and 10 discussed various technologies that can enable any processed social signals to be delivered to people who are blind and visually impaired, without overloading any of the their senses, like hearing or touch. Finally, Chapter 11 highlights the need for researchers to consider the impact of social mediation technologies on the society. While the discussions do not impart strict policies, this chapter initiates a conversation towards adopting important technology policies in the emerging assistive technology domains.

The dissertation presents a framework towards understanding the importance of social communication cues in various personal and professional environments. The dissertation presents an empirical evidence of the need for non-verbal cue enhancement for people who are blind and visually impaired. To this end, the technologies presented here only represent one possible approach towards addressing the important issues of disconnected social non-verbal communication. This dissertation establishes a framework for developing social mediation technologies each chapter from 4 through 11 presents possible future directions in each of the individual directions. While these technologies represent a

step towards advancing social mediation technologies, the material in this dissertation also presents a policy discussion towards developing and promoting wearable cameras within the everyday personal setting of our society. While a coarse level discussion of the policies are presented here, it lays a groundwork for the various considerations that are needed in propagating wearable cameras and social interpretation technologies into the future.

REFERENCES

- [1] N. Ambady and R. Rosenthal, "Thin slices of expressive behavior as predictors of interpersonal consequences : a Meta-Analysis," *Psychological Bulletin*, vol. 111, no. 2, pp. 274, 256, 1992.
- [2] A. Nakamura, T. Yamada, A. Goto, T. Kato, K. Ito, Y. Abe, T. Kachi, and R. Kakigi, "Somatosensory homunculus as drawn by MEG," *NeuroImage*, vol. 7, pp. 377–386, May 1998. PMID: 9626677.
- [3] R. E. Riggio, "Assessment of basic social skills," *Journal of Personality and Social Psychology*, vol. 51, no. 3, pp. 649–660, 1986.
- [4] B. D. Ruben, *Human communication handbook*. (Rochelle Park, N.J): Hayden Book Co., 1975.
- [5] M. L. Knapp and J. A. Hall, *Nonverbal Communication in Human Interaction*. Harcourt College Pub, 4th ed., Nov. 1996.
- [6] P. Borkenau, N. Mauer, R. Riemann, F. Spinath, and A. Angleitner, "Thin slices of behavior as cues of personality and intelligence.," *Journal of personality and social psychology*, vol. 86, no. 4, pp. 614, 599, 2004.
- [7] R. Brown, *Social Psychology*. New York, NY: Free Press, 1986.
- [8] J. Burgoon, D. Buller, J. Hale, and M. Turck, "Relational messages associated with nonverbal behaviors," *Human Communication Research*, vol. 10, no. 3, pp. 351–378, 1984.
- [9] C. Wetzel, "The midas touch: The effects of interpersonal touch on restaurant tipping," *Personality and Social Psychology Bulletin*, vol. 10, no. 4, pp. 512–517, 1984.
- [10] A. Haans and W. IJsselsteijn, "Mediated social touch: a review of current research and future directions," *Virtual Real.*, vol. 9, no. 2, pp. 149–159, 2006.
- [11] J. Bailenson and N. Yee, "Virtual interpersonal touch: Haptic interaction and copresence in collaborative virtual environments," *Multimedia Tools and Applications*, vol. 37, pp. 5–14, Mar. 2008.
- [12] A. J. Sameroff and M. J. Chandler, "Reproductive risk and the continuum of caretaker casualty," in *Review of Child Development Research* (F. D. Horowitz, ed.), vol. 4, Chicago: University of Chicago Press, 1975.

- [13] U. Altmann, R. Hermkes, and L. Alisch, "Analysis of nonverbal involvement in dyadic interactions," in *Verbal and Nonverbal Communication Behaviours*, pp. 37–50, 2007.
- [14] M. Zancanaro, B. Lepri, and F. Pianesi, "Automatic detection of group functional roles in face to face interactions," (Banff, Alberta, Canada), pp. 28–34, ACM, 2006.
- [15] W. Dong, B. Lepri, A. Cappelletti, A. S. Pentland, F. Pianesi, and M. Zancanaro, "Using the influence model to recognize functional roles in meetings," in *Proceedings of the 9th international conference on Multimodal interfaces*, (Nagoya, Aichi, Japan), pp. 271–278, ACM, 2007.
- [16] J. Hawkins and S. Blakeslee, *On Intelligence*. Times Books, adapted ed., Oct. 2004.
- [17] E. Rogers, W. Hart, and Y. Miike, "Edward t. hall and the history of intercultural communication: The united states and japan," *Keio Communication Review*, vol. 24, pp. 26, 3, 2002.
- [18] O. Hargie, *Social Skills in Interpersonal Communication*. Routledge, 3 ed., June 1994.
- [19] W. B. Walsh, K. H. Craik, and R. H. Price, *Person-environment psychology*. Routledge, 2000.
- [20] D. T. Kenrick and S. W. MacFarlane, "Ambient temperature and horn honking: A field study of the Heat/Aggression relationship," *Environment and Behavior*, vol. 18, pp. 179–191, Mar. 1986.
- [21] E. Krupat, *People in Cities: The Urban Environment and its Effects*. Cambridge University Press, Sept. 1985.
- [22] R. Sommer, *Personal Space: The Behavioral Basis of Design*. Prentice Hall Trade, 6th printing ed., June 1969.
- [23] R. Sommer, *Tight spaces; hard architecture and how to humanize it*. Prentice-Hall, 1974.
- [24] A. Schauss, "The psysiological effect of color on the suppression of human aggression," *International Journal of Biosocial Research*, vol. 7, pp. 55–64, 1985.
- [25] P. A. Bottomley and J. R. Doyle, "The interactive effects of colors and products on perceptions of brand logo appropriateness," *Marketing Theory*, vol. 6, pp. 63–83, Mar. 2006.

- [26] T. Farrenkopf and V. Roth, "The university faculty office as an environment.," *Environment and Behavior*, vol. 12, pp. 467–77, Dec. 1980.
- [27] R. H. Moos, *The Human Context: Environmental Determinants of Behavior*. Krieger Pub Co, June 1985.
- [28] V. Manusov and J. H. Harvey, *Attribution, Communication Behavior, and Close Relationships*. Cambridge University Press, 1 ed., Jan. 2001.
- [29] A. C. North, D. J. Hargreaves, and J. McKendrick, "In-store music affects product choice," *Nature*, vol. 390, p. 132, Nov. 1997.
- [30] J. Meer, "The light touch," *Psychology Today*, vol. 19, pp. 60–67, 1985.
- [31] D. S. Berry, "Attractive faces are not all created equal: Joint effects of facial babyishness and attractiveness on social perception," *Pers Soc Psychol Bull*, vol. 17, pp. 523–531, Oct. 1991.
- [32] B. H. Johnson, R. H. Nagasawa, and K. Peters, "Clothing style differences: Their effect on the impression of sociability," *Family and Consumer Sciences Research Journal*, vol. 6, pp. 58–63, Sept. 1977.
- [33] H. H. Jennings, *Sociometry in group relations*. (Washington): American Council on Education, 105 p. ed., 1959.
- [34] L. A. Zebrowitz, *Reading Faces*. Boulder CO: Westview Press, 1997.
- [35] D. S. Berry and L. Z. McArthur, "Perceiving character in faces: the impact of age-related craniofacial changes on social perception," *Psychological Bulletin*, vol. 100, pp. 3–18, July 1986. PMID: 3526376.
- [36] J. B. Corts and F. M. Gatti, "Physique and self-description of temperament," *Journal of Consulting Psychology*, vol. 29, pp. 432–439, Oct. 1965. PMID: 5827516.
- [37] L. A. Tucker, "Physical attractiveness, somatotype, and the male personality: A dynamic interactional perspective.," *Journal of Clinical Psychology*, vol. 40, no. 5, pp. 1226–34, 1984.
- [38] C. Cameron, S. Oskamp, and W. Sparks, "Courtship american style: Newspaper ads," *The Family Coordinator*, vol. 26, pp. 27–30, Jan. 1977. ArticleType: primary_article / Full publication date: Jan., 1977 / Copyright 1977 National Council on Family Relations.

- [39] C. L. Ogden, K. M. Flegal, M. D. Carroll, and C. L. Johnson, "Prevalence and trends in overweight among US children and adolescents, 1999-2000," *JAMA*, vol. 288, pp. 1728–1732, Oct. 2002.
- [40] J. H. Griffin, R. Bonazzi, J. H. Griffin, and R. Bonazzi, *Black Like Me*. Signet, 35th anniversary ed., Nov. 1996.
- [41] R. Porter, "Olfaction and human kin recognition," *Genetica*, vol. 104, pp. 259–263, Dec. 1998.
- [42] T. Lord and M. Kasprzak, "Identification of self through olfaction.," *Perceptual and motor skills*, vol. 69, no. 1, pp. 224, 219, 1989.
- [43] M. J. Russell, "Human olfactory communication," *Nature*, vol. 260, pp. 520–522, Apr. 1976.
- [44] N. Barber, "Mustache fashion covaries with a good marriage market for women," *Journal of Nonverbal Behavior*, vol. 25, pp. 261–272, Dec. 2001.
- [45] W. E. Hensley, "The effects of attire, location, and sex on aiding behavior: A similarity explanation," *Journal of Nonverbal Behavior*, vol. 6, no. 1, pp. 3–11, 1981.
- [46] N. Joseph, *Uniforms and Nonuniforms: Communication Through Clothing*. Greenwood Press, Nov. 1986.
- [47] T. L. Rosenfeld and T. G. Plax, "Clothing as communication," *Journal of Communication*, vol. 27, pp. 24–31.
- [48] C. Sanders and D. A. Vail, *Customizing the Body: The Art and Culture of Tattooing*. Temple University Press, Mar. 2008.
- [49] P. Ekman, "Nonverbal communication: Movements with precise meanings," 1976.
- [50] M. Wagner and N. Armstrong, *Field Guide to Gestures: How to Identify and Interpret Virtually Every Gesture Known to Man*. Quirk Books, July 2003.
- [51] D. Efron, *Gesture, Race and Culture*. Walter de Gruyter, Inc., Oct. 1972.
- [52] G. E. Weisfeld and J. M. Beresford, "Erectness of posture as an indicator of dominance or success in humans," *Motivation and Emotion*, vol. 6, pp. 113–131, June 1982.

- [53] E. C. Grant and J. H. Mackintosh, "A comparison of the social postures of some common laboratory rodents," *Behaviour*, vol. 21, no. 3/4, pp. 246–259, 1963. ArticleType: primary_article / Full publication date: 1963 / Copyright 1963 BRILL.
- [54] A. Kleinsmith, P. R. D. Silva, and N. Bianchi-Berthouze, "Cross-cultural differences in recognizing affect from body posture," *Interacting with Computers*, vol. 18, pp. 1371–1389, Dec. 2006.
- [55] A. Montagu, *Touching: The Human Significance of the Skin*. Harper Paperbacks, 3 ed., Sept. 1986.
- [56] W. A. Afifi and M. L. Johnson, "The use and interpretation of tie signs in a public setting: Relationship and sex differences," *Journal of Social and Personal Relationships*, vol. 16, pp. 9–38, Feb. 1999.
- [57] M. J. Hertenstein, J. M. Verkamp, A. M. Kerestes, and R. M. Holmes, "The communicative functions of touch in humans, nonhuman primates, and rats: a review and synthesis of the empirical research," *Genetic, Social, and General Psychology Monographs*, vol. 132, pp. 5–94, Feb. 2006. PMID: 17345871.
- [58] M. J. Hertenstein, D. Keltner, B. App, A. B. Bulleit, and R. Jaskolta, "Touch communicates distinct emotions," *Emotion*, vol. 6, no. 3, pp. 528–533, 2006.
- [59] G. Robles-De-La-Torre, "Principles of haptic perception in virtual environments," in *Human Haptic Perception: Basics and Applications*, pp. 363–379, 2008.
- [60] L. J. Carver and G. Dawson, "Development and neural bases of face recognition in autism," *Molecular Psychiatry*, vol. 7, no. s2, pp. S18–S20, 2002.
- [61] W. E. Rinn, "The neuropsychology of facial expression: A review of neurological and psychological mechanisms for producing facial expressions," *Psychological Bulletin*, vol. 95, pp. 52–77, 1984.
- [62] P. Ekman and W. Friesen, *Facial Action Coding System: A Technique for the Measurement of Facial Movement*. Consulting Psychologists Press, 1978.
- [63] C. E. Izard, *The maximally discriminative facial movement coding system*. Instructional Resources Center, University of Delaware, revised ed., 1983.
- [64] M. Argyle and M. Cook, *Gaze & Mutual Gaze*. Cambridge University Press, Jan. 1976.

- [65] C. L. Kleinke, "Gaze and eye contact: a research review," *Psychological Bulletin*, vol. 100, pp. 78–100, July 1986. PMID: 3526377.
- [66] A. Kendon, "Some functions of gaze-direction in social interaction.," *Acta Psychol (Amst)*, vol. 26, no. 1, pp. 63, 22, 1967.
- [67] M. S. Mast, "Dominance as expressed and inferred through speaking time," *Human Communication Research*, vol. 28, no. 3, pp. 420–450, 2002.
- [68] J. B. Bavelas, L. Coates, and T. Johnson, "Listener responses as a collaborative process: The role of gaze," *The Journal of Communication*, vol. 52, no. 3, pp. 566–580, 2002.
- [69] A. M. van Dulmen, P. F. M. Verhaak, and H. J. G. Bilo, "Shifts in Doctor-Patient communication during a series of outpatient consultations in Non-Insulin-Dependent diabetes mellitus.," *Patient Education and Counseling*, vol. 30, no. 3, pp. 227–37, 1997.
- [70] A. M. Glenberg, J. L. Schroeder, and D. A. Robertson, "Averting the gaze disengages the environment and facilitates remembering," *Memory & Cognition*, vol. 26, pp. 651–658, July 1998. PMID: 9701957.
- [71] J. Orozco, O. Rudovic, F. Roca, and J. Gonzalez, "Confidence assessment on eyelid and eyebrow expression recognition," in *Automatic Face & Gesture Recognition, 2008. FG '08. 8th IEEE International Conference on*, pp. 1–8, 2008.
- [72] C. Segrin and J. Flora, "Poor social skills are a vulnerability factor in the development of psychosocial problems.," *Human Communication Research*, vol. 26, no. 3, pp. 489–514, 2000.
- [73] D. Jindal-Snape, "Generalization and maintenance of social skills of children with visual impairments: Self-evaluation and the role of feedback," *Journal of Visual Impairment & Blindness*, vol. 98, pp. 470–483, 2004.
- [74] D. Jindal-Snape, "Use of feedback from sighted peers in promoting social interaction skills," *Journal of Visual Impairment and Blindness*, vol. 99, pp. 1–16, July 2005.
- [75] D. Jindal-Snape, "Using self-evaluation procedures to maintain social skills in a child who is blind," *Journal of Visual Impairment and Blindness*, vol. 92, pp. 362–366, 1998.
- [76] C. G. McGaha and D. C. Farran, "Interactions in an inclusive classroom: The effects of visual status and setting.," *Journal of Visual Impairment & Blindness*, vol. 95, pp. 80–94, 2001.

- [77] L. Kekelis, S. Sacks, and R. Gaylord-Ross, *The Development of Social Skills by Blind and Visually Impaired Students: Exploratory Studies and Strategies*. Amer Foundation for the Blind, June 1992.
- [78] K. Shinohara and J. Tenenberg, "A blind person's interactions with technology," *Commun. ACM*, vol. 52, no. 8, pp. 58–66, 2009.
- [79] K. Shinohara, "Designing assistive technology for blind users," in *Proceedings of the 8th international ACM SIGACCESS conference on Computers and accessibility*, (Portland, Oregon, USA), pp. 293–294, ACM, 2006.
- [80] K. Shinohara and J. Tenenberg, "Observing sara: a case study of a blind person's interactions with technology," in *Proceedings of the 9th international ACM SIGACCESS conference on Computers and accessibility*, (Tempe, Arizona, USA), pp. 171–178, ACM, 2007.
- [81] R. L. Daft and R. H. Lengel, "Organizational information requirements, media richness and structural design," *Manage. Sci.*, vol. 32, pp. 554–571, May 1986.
- [82] N. Kock and J. Nosek, "Expanding the boundaries of e-collaboration," *Professional Communication, IEEE Transactions on*, vol. 48, no. 1, pp. 1 – 9, 2005.
- [83] D. M. DeRosa, D. A. Hantula, N. Kock, and J. D'Arcy, "Trust and leadership in virtual teamwork: A media naturalness perspective," *Human Resource Management*, vol. 43, no. 2-3, pp. 219–232, 2004.
- [84] L. Robert and A. Dennis, "Paradox of richness: a cognitive model of media choice," *Professional Communication, IEEE Transactions on*, vol. 48, no. 1, pp. 10–21, 2005.
- [85] C. Solomon, "The challenges of working in virtual teams: Virtual teams survey report 2010," tech. rep., RW3 CultureWizard, New York, NY, 2010.
- [86] B. W. Tuckman, "Developmental sequence in small groups," *Psychological Bulletin*, vol. 63, pp. 384–399, 1965.
- [87] U. Hess and P. Philippot, *Group Dynamics and Emotional Expression*. Cambridge University Press, 1 ed., Jan. 2007.
- [88] N. Kock, "The ape that used email: Understanding e-communication behavior through evolution theory," *Communications of the Association for Information Systems*, vol. 5, no. 3, pp. 1–29, 2001.

- [89] V. L. Patel, J. Zhang, N. A. Yoskowitz, R. Green, and O. R. Sayan, “Translational cognition for decision support in critical care environments: a review,” *Journal of Biomedical Informatics*, vol. 41, pp. 413–431, June 2008. PMID: 18343731.
- [90] E. Salas, C. S. Burke, and K. C. Stagl, “Developing teams and team leaders: Strategies and principles.” in *Leader development for transforming organizations*. (R. G. Demaree, S. J. Zaccaro, and S. M. Halpin, eds.), pp. 325–358, Mahwah, NJ: Lawrence Erlbaum Associates, Inc., 2004.
- [91] P. Barach and M. Weingart, “Trauma team performance,” in *Trauma: Resuscitation, Anesthesia and Critical Care* (W. C. Wilson, C. M. Grande, and D. B. Hoyt, eds.), pp. 96–150, Informa Healthcare, 1 ed., Feb. 2007.
- [92] R. M. McIntyre and E. Salas, “Measuring and managing for team performance: Emerging principles from complex environments.” in *Team Effectiveness and Decision Making in Organizations* (R. A. Guzzo and E. Salas, eds.), pp. 194–203, San Francisco: Jossey-Bass, 1st ed., Mar. 1995.
- [93] J. A. Cannon-Bowers, S. I. Tannenbaum, and E. Salas, “Defining competencies and establishing team training requirements.” in *Team Effectiveness and Decision Making in Organizations* (R. A. Guzzo and E. Salas, eds.), pp. 333–380, San Francisco: Jossey-Bass, 1st ed., Mar. 1995.
- [94] J. E. Driskell and E. Salas, “Collective behavior and team performance,” *Hum. Factors*, vol. 34, no. 3, pp. 277–288, 1992.
- [95] D. Bandon, “Time to create sound teamwork,” *Journal of Qualitative Participation*, vol. 24, pp. 41–47, 2001.
- [96] H. King, J. Battles, D. Baker, A. Alonso, E. Salas, J. Webster, L. Toomey, and M. Salisbury, “TeamSTEPPS : Team strategies and tools to enhance performance and patient safety,” *Advances in Patient Safety: From Research to Implementation*, vol. 3.
- [97] A. S. Pentland, “Automatic mapping and modeling of human networks,” *PHYSICA A*, 2006.
- [98] T. Kim, D. O. Olguin, B. N. Waber, and A. S. Pentland, “Sensor-Based feedback systems in organizational computing,” in *Computational Science and Engineering, IEEE International Conference on*, vol. 4, (Los Alamitos, CA, USA), pp. 966–969, IEEE Computer Society, 2009.
- [99] L. Mann and A. Pirola-Merlo, “The relationship between individual creativity and team creativity: aggregating across people and time,” *Journal of Organizational Behavior*, vol. 25, no. 2, pp. 235–257, 2004.

- [100] J. Gardin, "Computer-Supported cooperative work: history and focus," *Computer*, vol. 27, no. 5, pp. 19–26, 1994.
- [101] G. B. Chapman and F. A. Sonnenberg, *Decision Making in Health Care: Theory, Psychology, and Applications*. Cambridge University Press, 1 ed., Sept. 2003.
- [102] S. W. J. Kozlowski, S. M. Gully, P. P. McHugh, E. Salas, and J. A. Cannon-Bowers, "A dynamic theory of leadership and team effectiveness: Developmental and task contingent leader roles.," in *Research in Personnel and Human Resources Management* (G. Ferris, ed.), vol. 19, pp. 252–305, JAI Press, 1 ed., 1996.
- [103] T. Sy, S. Ct, and R. Saavedra, "The contagious leader: impact of the leader's mood on the mood of group members, group affective tone, and group processes," *The Journal of Applied Psychology*, vol. 90, pp. 295–305, Mar. 2005. PMID: 15769239.
- [104] C. S.[1] and W. A., "Leadership of resuscitation teams: 'Lighthouse leadership'," *Resuscitation*, vol. 42, pp. 27–45, Sept. 1999.
- [105] S. Krishna, S. Panchanathan, and B. Patel, "Enriching interpersonal human interactions towards effective personal and professional communications," (Vancouver, British Columbia, Canada), Dec. 2010.
- [106] S. Krishna and B. Patel, "Studying individuals and groups Socio-Emotional artifacts in medical teams towards improved patient safety: A TeamSTEPPS approach," Tech. Rep. TR-10-009, Arizona State University, Tempe AZ, USA, July 2010.
- [107] R. R. Seethala, E. C. Esposito, and B. S. Abella, "Approaches to improving cardiac arrest resuscitation performance," *Current Opinion in Critical Care*, vol. 16, no. 3, pp. 196–202, 2010.
- [108] K. Karlgren, A. Dahlstrm, and S. Ponzer, "Design of an annotation tool to support simulation training of medical teams," in *Times of Convergence. Technologies Across Learning Contexts*, pp. 179–184, 2008.
- [109] D. N. Carbine, N. N. Finer, E. Knodel, and W. Rich, "Video recording as a means of evaluating neonatal resuscitation performance," *Pediatrics*, vol. 106, pp. 654–658, Oct. 2000. PMID: 11015505.
- [110] E. Bergs, F. Rutten, T. Tadros, P. Krijnen, and I. Schipper, "Communication during trauma resuscitation: do we know what is happening?," *Injury*, vol. 36, no. 8, pp. 905–911, 2005.
- [111] J. L. Moreno, *Sociometry, Experimental Method and the Science of Society*. Beacon House, Jan. 1983.

- [112] S. W. J. Kozlowski, S. M. Gully, P. P. McHugh, E. Salas, and Cannon-Bowers, “A dynamic theory of leadership and team effectiveness: Developmental and task contingent leader roles,” *Research in Personnel and Human Resources Management*, vol. 4, pp. 253–305, 1996.
- [113] N. Kock, D. A. Hantula, S. C. Hayne, G. Saad, P. M. Todd, and R. T. Watson, “Introduction to darwinian perspectives on electronic communication,” *Professional Communication, IEEE Transactions on*, vol. 51, no. 2, pp. 133–146, 2008.
- [114] S. Timmermans and A. Mauck, “The promises and pitfalls of Evidence-Based medicine,” *Health Aff*, vol. 24, pp. 18–28, Jan. 2005.
- [115] E. T. Hall, “A system for the notation of proxemic behavior,” *American Anthropologist*, vol. 65, pp. 1003–1026, Oct. 1963. ArticleType: primary_article / Issue Title: Selected Papers in Method and Technique / Full publication date: Oct., 1963 / Copyright 1963 American Anthropological Association.
- [116] Z. Zeng, M. Pantic, G. Roisman, and T. Huang, “A survey of affect recognition methods: Audio, visual, and spontaneous expressions,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 31, no. 1, pp. 39–58, 2009.
- [117] S. ur Rehman, L. Liu, and H. Li, “Manifold of facial expressions for tactile perception,” pp. 239–242, 2007.
- [118] A. Teeters, R. Kaliouby, and R. Picard, “Self-Cam: feedback from what would be your social partner,” in *SIGGRAPH '06: ACM SIGGRAPH 2006 Research posters*, (Boston, Massachusetts), p. 138, ACM, 2006.
- [119] A. Vinciarelli, M. Pantic, H. Bourlard, and A. Pentland, “Social signals, their function, and automatic analysis: a survey,” in *Proceedings of the 10th international conference on Multimodal interfaces*, (Chania, Crete, Greece), pp. 61–68, ACM, 2008.
- [120] T. Kim, A. Chang, L. Holland, and A. Pentland, “Meeting mediator: Enhancing group collaboration and leadership with sociometric feedback,” (San Diego, CA, USA), pp. 457–466, 2008.
- [121] A. Pentland, *Honest Signals: How They Shape Our World*. The MIT Press, Oct. 2008.
- [122] A. Vinciarelli, M. Pantic, H. Bourlard, and A. Pentland, “Social signal processing: state-of-the-art and future perspectives of an emerging domain,” in *Proceeding of the 16th ACM international conference on Multimedia*, (Vancouver, British Columbia, Canada), pp. 1061–1070, ACM, 2008.

- [123] R. E. Transon, "Using the feedback band device to control rocking behavior," *Journal of Visual Impairment & Blindness*, vol. 82, pp. 287 – 289, 1988.
- [124] J. N. Felps and R. J. Devlin, "Modification of stereotypic rocking of a blind adult.," *Journal of Visual Impairment and Blindness*, vol. 82, no. 3, pp. 107–08, 1988.
- [125] A. Adam, E. Rivlin, and I. Shimshoni, "Aggregated dynamic background modeling," pp. 3313–3316, 2006.
- [126] M. Yang, D. Kriegman, and N. Ahuja, "Detecting faces in images: a survey," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 24, no. 1, pp. 34–58, 2002.
- [127] P. Viola and M. Jones, "Robust real-time object detection," *International Journal of Computer Vision*, 2001.
- [128] R. Radke, S. Andra, O. Al-Kofahi, and B. Roysam, "Image change detection algorithms: a systematic survey," *Image Processing, IEEE Transactions on*, vol. 14, no. 3, pp. 294–307, 2005.
- [129] K. E. Ozden and L. V. Gool, "Background recognition in dynamic scenes with motion constraints," pp. 250–255, IEEE Computer Society, 2005.
- [130] Y. Ren, C. Chua, and Y. Ho, "Statistical background modeling for non-stationary camera," *Pattern Recogn. Lett.*, vol. 24, no. 1-3, pp. 183–196, 2003.
- [131] S. Todorovic and M. C. Nechyba, "Detection of artificial structures in Natural-Scene images using dynamic trees," pp. 35–39, IEEE Computer Society, 2004.
- [132] M. Barnard, M. Matilainen, and J. Heikkila, "Body part segmentation of noisy human silhouette images," pp. 1189–1192, 2008.
- [133] P. Srinivasan and J. Shi, "Bottom-up recognition and parsing of the human body," pp. 1–8, 2007.
- [134] A. Balan, L. Sigal, M. Black, J. Davis, and H. Houssecker, "Detailed human shape and pose from images," in *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, pp. 1–8, 2007.
- [135] X. Li, S. Maybank, S. Yan, D. Tao, and D. Xu, "Gait components and their application to gender recognition," *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, vol. 38, no. 2, pp. 145–155, 2008.

- [136] C. BenAbdelkader, R. Cutler, and L. Davis, "Person identification using automatic height and stride estimation," in *Pattern Recognition, 2002. Proceedings. 16th International Conference on*, vol. 4, pp. 377–380 vol.4, 2002.
- [137] R. Collins, R. Gross, and J. Shi, "Silhouette-based human identification from body shape and gait," in *Automatic Face and Gesture Recognition, 2002. Proceedings. Fifth IEEE International Conference on*, pp. 366–371, 2002.
- [138] A. Gallagher and T. Chen, "Clothing cosegmentation for recognizing people," pp. 1–8, 2008.
- [139] W. Zhang, B. Begole, M. Chu, J. Liu, and N. Yee, "Real-time clothes comparison based on multi-view vision," in *Distributed Smart Cameras, 2008. ICDS 2008. Second ACM/IEEE International Conference on*, pp. 1–10, 2008.
- [140] Y. Yacoob and L. Davis, "Detection and analysis of hair," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 28, no. 7, pp. 1164–1169, 2006.
- [141] Y. Xiao, N. Werghi, and P. Siebert, "A topological approach for segmenting human body shape," in *Image Analysis and Processing, 2003.Proceedings. 12th International Conference on*, pp. 82–87, 2003.
- [142] T. Sano and H. Yamamoto, "Human body shape imaging for japanese kimono design," vol. 2, pp. 1120–1123 Vol.2, 2004.
- [143] J. Lee, A. Jain, and R. Jin, "Scars, marks and tattoos (SMT): soft biometric for suspect and victim identification," in *Biometrics Symposium, 2008. BSYM '08*, pp. 1–8, 2008.
- [144] J. Yan and M. Pollefeys, "A Factorization-Based approach for articulated nonrigid shape, motion and kinematic chain recovery from video," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 30, no. 5, pp. 865–877, 2008.
- [145] J. Sung and D. Kim, "Combining local and global motion estimators for robust face tracking," pp. 345–350, 2007.
- [146] S. Du and R. Ward, "A robust approach for eye localization under variable illuminations," vol. 1, pp. I – 377–I – 380, 2007.
- [147] H. Lu and H. Lin, "Gender recognition using adaboosted feature," vol. 2, pp. 646–650, Third International Conference on Natural Computation, 2007. ICNC 2007., 2007.

- [148] H. Ling, S. Soatto, N. Ramanathan, and D. Jacobs, "A study of face recognition as people age," pp. 1–8, 2007.
- [149] X. Li, S. Maybank, and D. Tao, "Gender recognition based on local body motions," pp. 3881–3886, 2007.
- [150] F. Matta and J. Dugelay, "A behavioural approach to person recognition," pp. 1461–1464, IEEE International Conference on Multimedia and Expo, 2006, 2006.
- [151] M. Hahnel, D. Klunder, and K. Kraiss, "Color and texture features for person recognition," vol. 1, p. 652, 2004.
- [152] U. Saeed, F. Matta, and J. Dugelay, "Person recognition based on head and mouth dynamics," in *Multimedia Signal Processing, 2006 IEEE 8th Workshop on*, pp. 29–32, 2006.
- [153] U. Saeed and J. Dugelay, "Person recognition from video using facial mimics," in *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*, vol. 1, pp. I–493–I–496, 2007.
- [154] M. Kawade, "Vision-based face understanding technologies and applications," in *Micromechatronics and Human Science, 2002. MHS 2002. Proceedings of 2002 International Symposium on*, pp. 27–32, 2002.
- [155] A. Kanaujia, Y. Huang, and D. Metaxas, "Emblem detections by tracking facial features," in *Computer Vision and Pattern Recognition Workshop, 2006. CVPRW '06. Conference on*, p. 108, 2006.
- [156] H. Gunes and M. Piccardi, "Affect recognition from face and body: early fusion vs. late fusion," vol. 4, pp. 3443–3447, IEEE International Conference on Systems, Man and Cybernetics, 2005, 2005.
- [157] D. Kulic, W. Takano, and Y. Nakamura, "Combining automated on-line segmentation and incremental clustering for whole body motions," pp. 2591–2598, 2008.
- [158] N. Oliver, B. Rosario, and A. Pentland, "A bayesian computer vision system for modeling human interactions," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 22, no. 8, pp. 831–843, 2000.
- [159] J. Hwang, I. Karliga, and H. Cheng, "An automatic three-dimensional human behavior analysis system for video surveillance applications," p. 4 pp., IEEE International Symposium on Circuits and Systems, 2006. ISCAS 2006. Proceedings. 2006, 2006.

- [160] S. Park and J. Aggarwal, "Segmentation and tracking of interacting human body parts under occlusion and shadowing," in *Motion and Video Computing, 2002. Proceedings. Workshop on*, pp. 105–111, 2002.
- [161] S. Mitra and T. Acharya, "Gesture recognition: A survey," *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, vol. 37, no. 3, pp. 311–324, 2007.
- [162] J. Nunamaker, G. Tsechpenakis, D. Metaxas, M. Adkins, J. Kruse, J. Burgoon, M. Jensen, T. Meservy, D. Twitchell, and A. Deokar, "HMM-Based deception recognition from visual cues," pp. 824–827, *IEEE International Conference on Multimedia and Expo, 2005. ICME 2005.*, 2005.
- [163] I. Bacivarov, M. Ionita, and P. Corcoran, "Statistical models of appearance for eye tracking and eye-blink detection and measurement," *Consumer Electronics, IEEE Transactions on*, vol. 54, no. 3, pp. 1312–1320, 2008.
- [164] T. Meservy, M. Jensen, J. Kruse, J. Burgoon, J. Nunamaker, D. Twitchell, G. Tsechpenakis, and D. Metaxas, "Deception detection through automatic, unobtrusive analysis of nonverbal behavior," *Intelligent Systems, IEEE*, vol. 20, no. 5, pp. 36–43, 2005.
- [165] T. Funahashi, T. Fujiwara, and H. Koshimizu, "Face and eye tracking for gaze analysis," in *Control, Automation and Systems, 2007. ICCAS '07. International Conference on*, pp. 1337–1341, 2007.
- [166] A. Villanueva and R. Cabeza, "A novel gaze estimation system with one calibration point," *Systems, Man, and Cybernetics, Part B, IEEE Transactions on*, vol. 38, no. 4, pp. 1123–1138, 2008.
- [167] A. Fawky, S. Khalil, and M. Elsabrouty, "Eye detection to assist drowsy drivers," pp. 131–134, 2007.
- [168] U. Rajashekar, I. van der Linde, A. Bovik, and L. Cormack, "GAFFE: a Gaze-Attentive fixation finding engine," *Image Processing, IEEE Transactions on*, vol. 17, no. 4, pp. 564–573, 2008.
- [169] J. Rurainsky and P. Eisert, "Mirror-Based Multi-View analysis of facial motions," in *Image Processing, 2007. ICIIP 2007. IEEE International Conference on*, vol. 3, pp. III – 73–III – 76, 2007.
- [170] M. Yeasin, B. Bullot, and R. Sharma, "Recognition of facial expressions and measurement of levels of interest from video," *Multimedia, IEEE Transactions on*, vol. 8, no. 3, pp. 500–508, 2006.

- [171] S. Fukuda, "Detecting emotions and dangerous actions for better Human-System team working," pp. 205–206, Second International Conference on Secure System Integration and Reliability Improvement, 2008. SSIRI '08., 2008.
- [172] Y. Xie, Z. Wang, N. Cheng, G. Wang, and M. Nagai, "Facial and eye detection and application in affective recognition," in *Control Conference, 2006. CCC 2006. Chinese*, pp. 1942–1946, 2006.
- [173] T. J. Thompson, S. M. Pearcey, J. W. Bodfish, T. W. Crawford, and M. H. Lewis, "Stereotyped movement disorder in an adult following acquired brain injury: effect of environmental stimulation," *Behavioral Interventions*, vol. 10, pp. 79–85, Apr. 1995.
- [174] J. Jankovic, "Stereotypies," in *Movement Disorders* (C. D. Marsden and S. Fahn, eds.), vol. 3, pp. 503–517, London: Butterworth-Heinemann, 1994.
- [175] D. B. McAdam and C. M. O'Cleirigh, "Self-monitoring and verbal feedback to reduce stereotypic body rocking in a congenitally blind adult," *Re:View*, vol. 24, no. 4, p. 163, 1993.
- [176] R. S. Reivich and I. A. Rothrock, "Behavior problems of deaf children and adolescents: A Factor-Analytic study," *J Speech Hear Res*, vol. 15, pp. 93–104, Mar. 1972.
- [177] E. Haag, W. Huber, R. Hndgen, U. Stiller, and K. Willmes, "Repetitive verbal behavior in severe aphasia," *Der Nervenarzt*, vol. 56, pp. 543–52, Oct. 1985. PMID: 2415840.
- [178] D. B. Shabani, D. A. Wilder, and W. A. Flood, "Reducing stereotypic behavior through discrimination training, differential reinforcement of other behavior, and self-monitoring.," *Behavioral Interventions*, vol. 16, pp. 279–286, Oct. 2001.
- [179] R. L. Loftin, S. L. Odom, and J. F. Lantz, "Social interaction and repetitive motor behaviors.," *Journal of Autism & Developmental Disorders*, vol. 38, pp. 1124–1135, July 2008.
- [180] G. Yu, Y. Zhang, and R. Yan, "Loneliness, peer acceptance, and family functioning of chinese children with learning disabilities: Characteristics and relationships," *Psychology in the Schools*, vol. 42, no. 3, pp. 325–331, 2005.
- [181] V. J. Eichel, "A taxonomy for mannerisms of blind children.," *Journal of Visual Impairment and Blindness*, vol. 73, pp. 167–78, May 1979.
- [182] B. B. Blasch, "Blindisms: Treatment by punishment and reward in laboratory and natural settings," *Journal of Visual Impairment & Blindness*, pp. 215–230, 1972.

- [183] S. Raver, "Modification of head droop during conversation in a 3-Year-Old visually impaired child: A case study," *Journal of Visual Impairment and Blindness*, vol. 78, no. 7, pp. 307–10, 1984.
- [184] R. L. Ohlsen, "Control of body rocking in the blind through the use of vigorous exercise," *Journal of Instructional Psychology*, vol. 5, pp. 19–22, 1978.
- [185] A. H. Estevis and A. J. Koenig, "A cognitive approach to reducing stereotypic body rocking," *Re:View*, vol. 26, no. 3, p. 119, 1994.
- [186] P. J. Schloss and M. A. Smith, "Increasing appropriate behavior through related personal characteristics," in *Applied Behavior Analysis in the Classroom*, Boston: Allyn & Bacon, 1994.
- [187] G. Cartledge, *Teaching Social Skills to Children: Innovative Approaches*. Allyn & Bacon, 2 ed., June 1986.
- [188] S. Raver and P. W. Darsh, "Increasing social skills training for visually impaired children," *Education of the Visually Handicapped*, vol. 19, pp. 147–155, 1988.
- [189] S. Krishna, D. Colbry, J. Black, V. Balasubramanian, and S. Panchanathan, "A systematic requirements analysis and development of an assistive device to enhance the social interaction of people who are blind or visually impaired," in *Workshop on Computer Vision Applications for the Visually Impaired (CVAVI 08), European Conference on Computer Vision ECCV 2008*, (Marseille, France), Oct. 2008.
- [190] L. Bao and S. S. Intille, "Activity recognition from user-annotated acceleration data," pp. 1–17, 2004.
- [191] N. Ravi, N. Dandekar, P. Mysore, and M. Littman, "Activity recognition from accelerometer data," *American Association for Artificial Intelligence*, 2005.
- [192] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press, 1 ed., Mar. 2000.
- [193] L. Rabiner, "A tutorial on hidden markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [194] O. Amft, H. Junker, and G. Troster, "Detection of eating and drinking arm gestures using inertial body-worn sensors," pp. 160–163, 2005.

- [195] G. Chambers, S. Venkatesh, G. West, and H. Bui, “Hierarchical recognition of intentional human gestures for sports video annotation,” vol. 2, pp. 1082–1085 vol.2, 2002.
- [196] S. Lee and K. Mase, “Activity and location recognition using wearable sensors,” *Pervasive Computing, IEEE*, vol. 1, no. 3, pp. 24–32, 2002.
- [197] F. Foerster, M. Smeja, and J. Fahrenberg, “Detection of posture and motion by accelerometry: a validation in amulatory monitoring,” *Computer in Human Behavior*, vol. 15, pp. 571–583, 1999.
- [198] S. Arteaga, J. Chevalier, A. Coile, A. W. Hill, S. Sali, S. Sudhakhrisnan, and S. H. Kurniawan, “Low-cost accelerometry-based posture monitoring system for stroke survivors,” (Halifax, Nova Scotia, Canada), pp. 243–244, ACM, 2008.
- [199] N. Krishnan and S. Panchanathan, “Analysis of low resolution accelerometer data for continuous human activity recognition,” in *Acoustics, Speech and Signal Processing, 2008. ICASSP 2008. IEEE International Conference on*, pp. 3337–3340, 2008.
- [200] R. Polikar, “Ensemble based systems in decision making,” *Circuits and Systems Magazine, IEEE*, vol. 6, no. 3, pp. 21–45, 2006.
- [201] A. Vezhnevets and V. Vezhnevets, “Modest AdaBoost - teaching AdaBoost to generalize better,” (Novosibirsk Akademgorodok, Russia), 2005.
- [202] “DRM103 designer reference manual,” Tech. Rep. DRM103 Rev. 1, Freescale Semiconductor, Aug. 2008.
- [203] S. Krishna, G. Little, J. Black, and S. Panchanathan, “A wearable face recognition system for individuals with visual impairments,” in *Assets '05: Proceedings of the 7th international ACM SIGACCESS conference on Computers and accessibility*, (New York, NY, USA), pp. 106–113, ACM Press, 2005.
- [204] Y. Yuan and M. J. Shaw, “Induction of fuzzy decision trees,” *Fuzzy Sets Syst.*, vol. 69, no. 2, pp. 125–139, 1995.
- [205] Y. Benjamini, “Opening the box of a boxplot,” *The American Statistician*, vol. 42, pp. 257–262, Nov. 1988. ArticleType: primary_article / Full publication date: Nov., 1988 / Copyright 1988 American Statistical Association.
- [206] T. Fawcett, “An introduction to ROC analysis,” *Pattern Recogn. Lett.*, vol. 27, no. 8, pp. 861–874, 2006.

- [207] K. M. Newell, T. Incledon, and J. W. Bodfish, "Variability of stereotypic body-rocking in adults with mental retardation," *American Journal on Mental Retardation*, vol. 104, pp. 279 – 88, May 1999.
- [208] R. Raina, A. Battle, H. Lee, B. Packer, and A. Y. Ng, "Self-taught learning: transfer learning from unlabeled data," in *Proceedings of the 24th international conference on Machine learning*, (Corvalis, Oregon), pp. 759–766, ACM, 2007.
- [209] S. Wang, Z. Liu, S. Lv, Y. Lv, G. Wu, P. Peng, F. Chen, and X. Wang, "A natural visible and infrared facial expression database for expression recognition and emotion inference," *Multimedia, IEEE Transactions on*, vol. 12, no. 7, pp. 682 –691, 2010.
- [210] G. Jiang, X. Song, F. Zheng, P. Wang, and A. Omer, "Facial expression recognition using thermal image," *Conference Proceedings: ... Annual International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conference*, vol. 1, pp. 631–633, 2005. PMID: 17282261.
- [211] G. Gibert, M. Pruzinec, T. Schultz, and C. Stevens, "Enhancement of human computer interaction with facial electromyographic sensors," in *Proceedings of the 21st Annual Conference of the Australian Computer-Human Interaction Special Interest Group: Design: Open 24/7, OZCHI '09*, (New York, NY, USA), pp. 421–424, ACM, 2009.
- [212] P. Ekman, "An argument for basic emotions," *Cognition & Emotion*, vol. 6, no. 3, pp. 169–200, 1992.
- [213] Z. Zeng, M. Pantic, G. I. Roisman, and T. S. Huang, "A survey of affect recognition methods: Audio, visual, and spontaneous expressions," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 1, pp. 39–58, 2009.
- [214] A. B. Ashraf, S. Lucey, J. F. Cohn, T. Chen, Z. Ambadar, K. M. Prkachin, and P. E. Solomon, "The painful face - pain expression recognition using active appearance models," *Image Vision Comput.*, vol. 27, no. 12, pp. 1788–1796, 2009.
- [215] M. Bartlett, G. Littlewort, P. Braathen, T. Sejnowski, and J. Movellan, "A prototype for automatic recognition of spontaneous facial actions"," vol. 15, pp. 1271–1278, 2003.
- [216] M. Yeasin, B. Bullot, and R. Sharma, "Recognition of facial expressions and measurement of levels of interest from video," *Multimedia, IEEE Transactions on*, vol. 8, no. 3, pp. 500–508, 2006.

- [217] M. S. Bartlett, G. Littlewort, M. Frank, C. Lainscsek, I. Fasel, and J. Movellan, "Fully automatic facial action recognition in spontaneous behavior," pp. 223–230, IEEE Computer Society, 2006.
- [218] I. Cohen, N. Sebe, A. Garg, L. S. Chen, and T. S. Huang, "Facial expression recognition from video sequences: temporal and static modeling," *Computer Vision and Image Understanding*, vol. 91, pp. 160–187, Aug. 2003.
- [219] J. Cohn, L. Reed, Z. Ambadar, J. Xiao, and T. Moriyama, "Automatic analysis and recognition of brow actions and head motion in spontaneous facial behavior," vol. 1, pp. 610–616 vol.1, 2004.
- [220] R. E. Kaliouby and P. Robinson, "Real-Time inference of complex mental states from facial expressions and head gestures," p. 154, IEEE Computer Society, 2004.
- [221] B. Fasel, F. Monay, and D. Gatica-Perez, "Latent semantic analysis of facial action codes for automatic facial expression recognition," (New York, NY, USA), pp. 181–188, ACM, 2004.
- [222] S. V. Ioannou, A. T. Raouzaïou, V. A. Tzouvaras, T. P. Mailis, K. C. Karpouzis, and S. D. Kollias, "Emotion recognition through facial expression analysis based on a neurofuzzy network," *Neural Netw.*, vol. 18, no. 4, pp. 423–435, 2005.
- [223] Q. Ji, P. Lan, and C. Looney, "A probabilistic framework for modeling and real-time monitoring human fatigue," *Systems, Man and Cybernetics, Part A: Systems and Humans, IEEE Transactions on*, vol. 36, no. 5, pp. 862–875, 2006.
- [224] A. Kapoor and R. W. Picard, "Multimodal affect recognition in learning environments," (Hilton, Singapore), pp. 677–682, ACM, 2005.
- [225] A. Kapoor, W. Bursleson, and R. W. Picard, "Automatic prediction of frustration," *Int. J. Hum.-Comput. Stud.*, vol. 65, no. 8, pp. 724–736, 2007.
- [226] C. Lee and A. Elgammal, "Facial expression analysis using nonlinear decomposable generative models," in *Analysis and Modelling of Faces and Gestures*, pp. 17–31, 2005.
- [227] G. C. Littlewort, M. S. Bartlett, and K. Lee, "Faces of pain: automated measurement of spontaneous facial expressions of genuine and posed pain," (Nagoya, Aichi, Japan), pp. 15–21, ACM, 2007.
- [228] S. Lucey, A. Bilal, and J. Cohn, "Investigating spontaneous facial action recognition through AAM representations of the face," in *Face Recognition Book* (K. Kurihara, ed.), Pro Literatur Verlag, 2007.

- [229] M. Pantic and I. Patras, "Dynamics of facial expression: recognition of facial actions and their temporal segments from face profile image sequences," *IEEE Transactions on Systems, Man, and Cybernetics. Part B, Cybernetics: A Publication of the IEEE Systems, Man, and Cybernetics Society*, vol. 36, pp. 433–449, Apr. 2006. PMID: 16602602.
- [230] M. Pantic and L. Rothkrantz, "Case-based reasoning for user-profiled recognition of emotions from face images," vol. 1, pp. 391–394 Vol.1, 2004.
- [231] N. Sebe, M. Lew, I. Cohen, Y. Sun, T. Gevers, and T. Huang, "Authentic facial expression analysis," pp. 517–522, 2004.
- [232] Y. Tong, W. Liao, and Q. Ji, "Facial action unit recognition by exploiting their dynamic and semantic relationships," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 29, no. 10, pp. 1683–1699, 2007.
- [233] M. F. Valstar, H. Gunes, and M. Pantic, "How to distinguish posed from spontaneous smiles using geometric features," (Nagoya, Aichi, Japan), pp. 38–45, ACM, 2007.
- [234] M. F. Valstar, M. Pantic, Z. Ambadar, and J. F. Cohn, "Spontaneous vs. posed facial behavior: automatic analysis of brow actions," (Banff, Alberta, Canada), pp. 162–170, ACM, 2006.
- [235] H. Wang and N. Ahuja, "Facial expression decomposition," p. 958, IEEE Computer Society, 2003.
- [236] J. Wang, L. Yin, X. Wei, and Y. Sun, "3D facial expression recognition based on primitive surface feature distribution," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Volume 2 (CVPR'06)*, (New York, NY, USA), pp. 1399–1406, 2006.
- [237] Z. Wen and T. Huang, "Capturing subtle facial motions in 3D face tracking," pp. 1343–1350 vol.2, 2003.
- [238] J. Whitehill and C. W. Omlin, "Haar features for FACS AU recognition," pp. 97–101, IEEE Computer Society, 2006.
- [239] Z. Zeng, Y. Fu, G. Roisman, Z. Wen, Y. Hu, and T. Huang, "Spontaneous emotional facial expression detection," *Journal of Multimedia*, vol. 1, no. 5, pp. 1–8, 2006.
- [240] S. Krishna, T. McDaniel, and S. Panchanathan, "Embodied social interaction assistant," tech. rep., Arizona State University, Tempe, USA, Jan. 2010.

- [241] L. Shang and K. Chan, "Temporal Exemplar-Based bayesian networks for facial expression recognition," in *Machine Learning and Applications, 2008. ICMLA '08. Seventh International Conference on*, pp. 16–22, 2008.
- [242] R. Paget, I. D. Longstaff, and B. Lovell, "Texture classification using nonparametric markov random fields," vol. 1, pp. 67–70, 1997.
- [243] Y. Tian, T. Kanade, and J. Cohn, "Recognizing action units for facial expression analysis," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 23, no. 2, pp. 97–115, 2001.
- [244] M. M. Sohlberg and C. A. Mateer, *Introduction to Cognitive Rehabilitation: Theory and Practice*. The Guilford Press, Feb. 1989.
- [245] M. Grunwald, *Human Haptic Perception: Basics and Applications*. Birkhuser Basel, 1 ed., Nov. 2008.
- [246] D. Maynes-aminzade, "Edible bits: Seamless interfaces between people, data and food," *Proceedings of the 2005 ACM Conference on Human Factors in Computing Systems (CHI'2005) - Invited Talk*, 2005.
- [247] Y. Chen, "Olfactory display: Development and application in virtual reality therapy," in *International Conference on Artificial Reality and Telexistence*, (Los Alamitos, CA, USA), pp. 580–584, IEEE Computer Society, 2006.
- [248] N. J. Dominy and P. W. Lucas, "Ecological importance of trichromatic vision to primates," *Nature*, vol. 410, pp. 363–366, Mar. 2001.
- [249] K. Koch, J. McLean, R. Segev, M. A. Freed, M. J. B. II, V. Balasubramanian, and P. Sterling, "How much the eye tells the brain," *Current Biology*, vol. 16, no. 14, pp. 1428–1434, 2006.
- [250] P. B. y rita, C. C. Collins, F. A. Saunders, B. White, and L. Scaden, "Vision substitution by tactile image projection," *Nature*, vol. 221, pp. 963–964, Mar. 1969.
- [251] S. Rehman, L. Liu, and H. Li, "Vibrotactile rendering of human emotions on the manifold of facial expressions," *Journal of Multimedia*, vol. 3, no. 3, pp. 18–25, 2008.
- [252] T. McDaniel, S. Krishna, V. Balasubramanian, D. Colbry, and S. Panchanathan, "Using a haptic belt to convey non-verbal communication cues during social interactions to individuals who are blind," in *Haptic Audio visual Environments and Games, 2008. HAVE 2008. IEEE International Workshop on*, pp. 13–18, 2008.

- [253] S. Brave and A. Dahley, “inTouch: a medium for haptic interpersonal communication,” in *CHI '97 extended abstracts on Human factors in computing systems: looking to the future*, (Atlanta, Georgia), pp. 363–364, ACM, 1997.
- [254] C. Dodge, “The bed: a medium for intimate communication,” in *CHI '97 extended abstracts on Human factors in computing systems: looking to the future*, (Atlanta, Georgia), pp. 371–372, ACM, 1997.
- [255] B. J. Fogg, L. D. Cutler, P. Arnold, and C. Eisbach, “HandJive: a device for interpersonal haptic entertainment,” in *Proceedings of the SIGCHI conference on Human factors in computing systems*, (Los Angeles, California, United States), pp. 57–64, ACM Press/Addison-Wesley Publishing Co., 1998.
- [256] O. Morikawa, J. Yamashita, and Y. Fukui, “The sense of physically crossing paths: creating a soft initiation in HyperMirror communication,” in *CHI '00 extended abstracts on Human factors in computing systems*, (The Hague, The Netherlands), pp. 183–184, ACM, 2000.
- [257] A. Chang, S. O’Modhrain, R. Jacob, E. Gunther, and H. Ishii, “ComTouch: design of a vibrotactile communication device,” in *DIS '02: Proceedings of the conference on Designing interactive systems*, pp. 320, 312, ACM Press, 2002.
- [258] A. Rovers and H. van Essen, “HIM: a framework for haptic instant messaging,” in *CHI '04 extended abstracts on Human factors in computing systems*, (Vienna, Austria), pp. 1313–1316, ACM, 2004.
- [259] F. F. Mueller, F. Vetere, M. R. Gibbs, J. Kjeldskov, S. Pedell, and S. Howard, “Hug over a distance,” in *CHI '05 extended abstracts on Human factors in computing systems*, (Portland, OR, USA), pp. 1673–1676, ACM, 2005.
- [260] K. Dobson, danah boyd, W. Ju, J. Donath, and H. Ishii, “Creating visceral personal and social interactions in mediated spaces,” in *CHI '01 extended abstracts on Human factors in computing systems*, (Seattle, Washington), pp. 151–152, ACM, 2001.
- [261] M. O. Alhalabi, S. Horiguchi, and S. Kunifuji, “An experimental study on the effects of network delay in cooperative shared haptic virtual environment,” *Computers & Graphics*, vol. 27, pp. 205–213, Apr. 2003.
- [262] L. Bonanni, C. Vaucelle, J. Lieberman, and O. Zuckerman, “TapTap: a haptic wearable for asynchronous distributed touch therapy,” in *CHI '06: CHI '06 extended abstracts on Human factors in computing systems*, (Montr\`eal, Qu\`ebec, Canada), pp. 585, 580, ACM, 2006.

- [263] J. Werner, R. Wettach, and E. Hornecker, “United-pulse: feeling your partner’s pulse,” in *Proceedings of the 10th international conference on Human computer interaction with mobile devices and services*, (Amsterdam, The Netherlands), pp. 535–538, ACM, 2008.
- [264] L. Cappelletti, M. Feeri, and G. Nicoletti, “Vibrotactile colour rendering for the visually impaired within the VIDET project,” in *Telemanipulator and Telepresence Technologies V*, vol. 3524, pp. 92–96, Nov. 1998.
- [265] T. Oron-Gilad, J. Downs, R. Gilson, and P. Hancock, “Vibrotactile guidance cues for target acquisition,” *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, vol. 37, no. 5, pp. 993–1004, 2007.
- [266] H. Uchiyama, M. A. Covington, and W. D. Potter, “Vibrotactile glove guidance for semi-autonomous wheelchair operations,” in *Proceedings of the 46th Annual Southeast Regional Conference on XX*, (Auburn, Alabama), pp. 336–339, ACM, 2008.
- [267] A. Hein and M. Brell, “conTACT - a vibrotactile display for computer aided surgery,” in *World Haptics Conference*, vol. 0, (Los Alamitos, CA, USA), pp. 531–536, IEEE Computer Society, 2007.
- [268] M. Brell, D. Rokamp, and A. Hein, “Fusion of vibrotactile signals used in a tactile display in computer aided surgery,” in *Haptics: Perception, Devices and Scenarios*, pp. 383–388, 2008.
- [269] J. Cha, Y. Seo, Y. Kim, and J. Ryu, “An Authoring/Editing framework for haptic broadcasting: Passive haptic interactions using MPEG-4 BIFS,” in *EuroHaptics Conference, 2007 and Symposium on Haptic Interfaces for Virtual Environment and Teleoperator Systems. World Haptics 2007. Second Joint*, pp. 274–279, 2007.
- [270] J. Cha, Y. Ho, Y. Kim, J. Ryu, and I. Oakley, “A framework for haptic broadcasting,” *IEEE MultiMedia*, vol. 16, no. 3, pp. 16–27, 2009.
- [271] A. M. Murray, R. L. Klatzky, and P. K. Khosla, “Psychophysical characterization and testbed validation of a wearable vibrotactile glove for telemanipulation,” *Presence: Teleoper. Virtual Environ.*, vol. 12, no. 2, pp. 156–182, 2003.
- [272] A. Schwaninger, C. C. Carbon, and H. Leder, “Expert face processing: Specilization and constraints,” 2003.
- [273] V. Bruce and A. Young, *In the Eye of the Beholder: The Science of Face Perception*. Oxford University Press, 2006.

- [274] J. Sadr, I. Jarudi, and P. Shinha, "The role of eyebrows in face recognition," *Perception*, vol. 32, pp. 285–93, 2003.
- [275] A. W. Young, D. Hellawell, and D. C. Hay, "Configurational information in face perception," *Perception*, vol. 16, pp. 747–59, 1987.
- [276] P. Sinha, B. J. Balas, Y. Ostrovsky, and R. Russell, "Face recognition by humans: 19 results all computer vision researchers should know about," *Proceedings of the IEEE*, vol. 94, no. 11, pp. 1948–62, 2006.
- [277] H. McConachie, "Developmental prosopagnosia: A single case report," *Cortex*, vol. 12, pp. 76–82, 1976.
- [278] I. Kennerknecht, T. Gruter, B. Welling, S. Wentzek, J. Horst, S. Edwards, and M. Gruter, "First report of prevalence of non-syndromic hereditary prosopagnosia (hpa)," *American Journal of Medical Genetics, Part A*, vol. 140, pp. 1617–22, 2006.
- [279] D. J. Turk, T. C. Handy, and M. S. Gazzaniga, "Can perceptual expertise account for the own-race bias in face recognition? a split brain study," *Cognitive Neuropsychology*, vol. 22, no. 7, pp. 877–83, 2005.
- [280] B. Gkberk, M. O. Irfanoglu, L. Akarun, and E. Alpaydin, "Learning the best subset of local features for face recognition," *Pattern Recognition*, vol. 40, no. 5, p. 1520, 2007.
- [281] L. Wiskott, "Phantom faces for face analysis," *Pattern Recognition*, vol. 30, no. 6, p. 837, 1997.
- [282] N. Kruger, M. Potzsch, and C. Malsburg, "Determination of face position and pose with a learned representation based on labelled graphs," *Image Vision Computing*, vol. 15, p. 665, 1997.
- [283] P. Kalocsai, C. Malsburg, and J. Horn, "Face recognition by statistical analysis of feature detectors," *Image Vision Computing*, vol. 18, no. 4, p. 273, 2000.
- [284] A. Tefas, C. Kotropoulos, and I. Pitas, "Using support vector machines to enhance the performance of elastic graph matching for frontal face authentication," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 7, p. 735, 2001.
- [285] D. H. Liu, K. M. Lam, and L. S. Shen, "Optimal sampling of gabor features for face recognition," *Pattern Recognition Letters*, vol. 25, no. 2, p. 267, 2004.

- [286] P. Yang, S. Shan, W. Gao, S. Li, and D. Zang, "Face recognition using ada-boosted gabor features," in *Proceedings of the 16th International Conference on Face and Gesture Recognition*, 2004.
- [287] X. Wang and H. Oi, "Face recognition using optimal non-orthogonal wavelet basis evaluated by information complexity," in *Proceedings of the 16th International Conference on Pattern Recognition*, p. 164, 2002.
- [288] L. Wiskott, J.-M. Fellous, N. Krüger, and C. von der Malsburg, "Face recognition by elastic bunch graph matching," (Heidelberg), pp. 456–463, Springer-Verlag, 1997.
- [289] C. von der Malsburg, "Nervous structures with dynamical links," *Ber. Bunsenges. Phys. Chem.*, vol. 89, pp. 703–710, 1985.
- [290] E. Bienenstock and christoph von der Malsburg, "A neural network for invariant pattern recognition," *Europhysics Letters*, vol. 4, pp. 121–126, 1987.
- [291] L. Wiskott, J.-M. Fellous, N. Krüger, and C. von der Malsburg, *Face Recognition and Gender Determination*. 1995. Published: International Workshop on Automatic Face- and Gesture-Recognition, Zürich, June 26-28, 1995.
- [292] L. Wiskott and C. von der Malsburg, *Face Recognition by Dynamic Link Matching*. <http://www.cs.utexas.edu/users/nn/web-pubs/htmlbook96/>: The UTCS Neural Networks Research Group, Austin, TX, 1996.
- [293] M. Lyons, M. Lyons, J. Budynek, J. Budynek, and S. Akamatsu, "Automatic classification of single facial images," *Transactions on Pattern Analysis and Machine Intelligence*, vol. 21, pp. 1357–1362, 1999.
- [294] Y. Liu and Chongqing, "Face recognition using kernel principal component analysis and genetic algorithms," in *12th IEEE Workshop on Neural Networks for Signal Processing*, p. 337, Sep 2002 2002.
- [295] Y. Xu, B. Li, and B. Wang, "Face recognition by fast independent component analysis and genetic algorithm," in *Fourth International Conference on Computer and Information*, p. 194, 14-16 Sep 2004 2004.
- [296] K. Wong and K. Lam, "A reliable approach for human face detection using genetic algorithm," in *IEEE International Symposium on Circuits and Systems*, vol. 4, p. 499, 1999.
- [297] S. Karungaru, M. Fukumi, and N. Akamatsu, "Face recognition using genetic algorithm based template matching," in *International Symposium on Communications and Information Technologies*, October 26- 29,2004 2004.

- [298] D. Ozkan, "Feature selection for face recognition using a genetic algorithm," tech. rep., 2006.
- [299] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 20, pp. 91–110, 2003.
- [300] J. Huang and H. Wechsler, "Eye location using genetic algorithm," in *2nd International Conference on Audio and Video-Based Biometric Person Authentication*, 1999.
- [301] Y. Sun and L. Yin, "A genetic algorithm based feature selection approach for 3d face recognition," in *Biometrics Consortium Conference*, September 19-21, 2005 2005.
- [302] Z. Sun, X. Yuan, G. Bebis, and S. Louis, "Neural-network-based gender classification using genetic eigen-feature extraction," 2002.
- [303] J. Black, M. Gargesha, K. Kahol, and S. Panchanathan, "A framework for performance evaluation of face recognition algorithms," *ITCOM, Internet Multimedia Systems II, Boston*, July 2002.
- [304] G. Little, S. Krishna, J. Black, and S. Panchanathan, "A methodology for evaluating robustness of face recognition algorithms with respect to variations in pose and illumination angle," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, (Philadelphia, USA), pp. 89–92, 2005.
- [305] L. Sirovich and M. Kirby, "Low-dimensional procedure for the characterization of human faces," *J. Opt. Soc. Am. A*, vol. 4, p. 519, 1987.
- [306] M. Turk and A. Pentland, "Face recognition using eigenfaces," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pp. 586–591, 1991.
- [307] E. T. Hall, *The Hidden Dimension*. Anchor, Oct. 1990.
- [308] P. Viola and M. J. Jones, "Robust Real-Time face detection," *Int. J. Comput. Vision*, vol. 57, no. 2, pp. 137–154, 2004.
- [309] E. Hjelm and B. K. Low, "Face detection: A survey," *Computer Vision and Image Understanding*, vol. 83, pp. 236–274, Sept. 2001.
- [310] A. Hadid and M. Pietikainen, "A hybrid approach to face detection under unconstrained environments," *18th International Conference on Pattern Recognition*, vol. 1, pp. 227–230, 2006.

- [311] I. Naseem and M. Deriche, "Robust human face detection in complex color images," *IEEE International Conference on Image Processing*, vol. 2, pp. 338–41, 2005.
- [312] M. Wimmer, B. Radig, and M. Beetz, "A person and context specific approach for skin color classification," *Pattern Recognition, 2006. ICPR 2006. 18th International Conference on*, vol. 2, pp. 39–42, 2006.
- [313] M. B. Hmid and Y. B. Jemaa, "Fuzzy classification, image segmentation and shape analysis for human face detection," *8th International Conference on Signal Processing*, vol. 4, 2006.
- [314] U. Tariq, H. Jamal, M. Shahid, and M. Malik, "Face detection in color images, a robust and fast statistical approach," *Proceedings of INMIC 2004. 8th International Multitopic Conference*, pp. 73–78, 2004.
- [315] Y.-W. Wu and X.-Y. Ai, "Face detection in color images using adaboost algorithm based on skin color information," *International Workshop on Knowledge Discovery and Data Mining*, pp. 339–342, 2008.
- [316] K. Sentz and S. Ferson, "Combination of evidence in dempster-shafer theory," tech. rep., Sandia National Laboratories, 2002.
- [317] P. J. Phillips, H. Moon, P. Rauss, and S. A. Rizvi, "The feret evaluation methodology for face-recognition algorithms," *Proceedings of the 1997 Conference on Computer Vision and Pattern Recognition*, p. 137, 1997.
- [318] J. Bilmes, "A gentle tutorial of the em algorithm and its application to parameter estimation for gaussian mixture and hidden markov models," (Berkeley CA), International Computer Science Institute, U.C. Berkeley, April, 1998.
- [319] P. Perez, "Markov random fields and images," *CWI Quaterly*, vol. 11, pp. 413–437, 1998.
- [320] M. Vezjak and M. Stephanic, "An anthropological model for automatic recognition of the male human face," *Annals of Human Biology*, vol. 21, pp. 363–380, 1994.
- [321] B. Leibe, A. Leonardis, and B. Schiele, "Combined object categorization and segmentation with an implicit shape model," *In ECCV Workshop on Statistical Learning in Computer Vision*, pp. 17–32, 2004.
- [322] F. Porikli and O. Tuzel, "Object tracking in Low-Frame-Rate video," *SPIE Image and Video Communications and Processing*, vol. 5685, pp. 72–79, Mar. 2005.

- [323] Y. Li, H. Ai, T. Yamashita, S. Lao, and M. Kawade, "Tracking in low frame rate video: A cascade particle filter with discriminative observers of different lifespans," in *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, pp. 1–8, 2007.
- [324] J. Kwon and K. M. Lee, "Tracking of abrupt motion using wang-landau monte carlo estimation," in *Proceedings of the 10th European Conference on Computer Vision: Part I, ECCV '08*, (Berlin, Heidelberg), pp. 387–400, Springer-Verlag, 2008.
- [325] P. Viola and M. Jones, "Robust real-time object detection," *International Journal of Computer Vision - to appear*, 2002.
- [326] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," *IN CVPR*, vol. 1, pp. 886–893, 2005.
- [327] K. Nummiaro, E. Koller-Meier, and L. V. Gool, "An adaptive color-based particle filter," *Image and Vision Computing*, vol. 21, pp. 99–110, Jan. 2003.
- [328] F. Porikli, "Integral histogram: A fast way to extract histograms in cartesian spaces," *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, vol. 1, pp. 829–836, 2005.
- [329] Q. Zhu, Q. Zhu, S. Avidan, S. Avidan, M. chen Yeh, M. chen Yeh, K. ting Cheng, and K. ting Cheng, "Fast human detection using a cascade of histograms of oriented gradients," *IN CVPR06*, vol. 2006, pp. 1491–1498, 2006.
- [330] F. Porikli, P. Meer, O. Tuzel, and O. Tuzel, "P.: Human detection via classification on riemannian manifolds," *In Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 1–8, 2007.
- [331] R. Fergus, P. Perona, and A. Zisserman, "Object class recognition by unsupervised scale-invariant learning," in *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, vol. 2, pp. II–264–II–271 vol.2, 2003.
- [332] C. Yang, R. Duraiswami, and L. Davis, "Fast multiple object tracking via a hierarchical particle filter," *In Proc. of Intl. Conf. on Computer Vision*, vol. 1, pp. 212–219, 2005.
- [333] H. G. Barrow, J. M. Tenenbaum, R. C. Bolles, and H. C. Wolf, "Parametric correspondence and chamfer matching: Two new techniques for image matching," 2006.
- [334] B. Leibe, E. Seemann, and B. Schiele, "Pedestrian detection in crowded scenes," *IN CVPR*, vol. 1, pp. 878–885, 2005.

- [335] M. S. Arulampalam, S. Maskell, and N. Gordon, "A tutorial on particle filters for on-line nonlinear/non-Gaussian bayesian tracking," *IEEE Trans. on Signal Processing*, vol. 50, pp. 174—188, 2002.
- [336] M. Isard and A. Blake, "CONDENSATION - conditional density propagation for visual tracking," *Intl. Journal of Computer Vision*, vol. 29, pp. 5—28, 1998.
- [337] S. Birchfield, "Elliptical head tracking using intensity gradients and color histograms," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 232—237, 1998.
- [338] K. Okuma, A. Taleghani, N. D. Freitas, O. D. Freitas, J. J. Little, and D. G. Lowe, "A boosted particle filter: Multitarget detection and tracking," *In Proc. of European Conf. on Computer Vision*, vol. 1, pp. 28—39, 2004.
- [339] F. Crow, "Summed-area tables for texture mapping," in *SIGGRAPH '84: Proceedings of the 11th annual conference on Computer graphics and interactive techniques*, pp. 207–212, ACM, 1984.
- [340] V. Manohar, P. Soundararajan, D. Goldgof, R. Kasturi, and J. Garofolo, "Performance evaluation of object detection and tracking in video," *In Proc. of Seventh Asian Conf. on Computer Vision*, vol. 2, pp. 151—161, 2006.
- [341] S. Brewster and L. Brown, "Tactons: structured tactile messages for non-visual information display," in *AUIC '04: Proceedings of the fifth conference on Australasian user interface*, pp. 23, 15, Australian Computer Society, Inc., 2004.
- [342] L. A. Jones and N. B. Sarter, "Tactile displays: guidance for their design and application," *Human Factors*, vol. 50, pp. 90–111, Feb. 2008. PMID: 18354974.
- [343] H. Z. Tan and A. Pentland, "Tactual displays for sensory substitution and wearable computers," in *ACM SIGGRAPH 2005 Courses on - SIGGRAPH '05*, (Los Angeles, California), p. 105, 2005.
- [344] F. A. Geldard, "Adventures in tactile literacy.," *American Psychologist*. Vol. 12(3), vol. 12, pp. 115–124, Mar. 1957.
- [345] L. A. Jones and K. Ray, "Localization and pattern recognition with tactile displays," in *Proceedings of the 2008 Symposium on Haptic Interfaces for Virtual Environment and Teleoperator Systems*, pp. 33–39, IEEE Computer Society, 2008.
- [346] F. A. Geldard and C. E. Sherrick, "The cutaneous "Rabbit": a perceptual illusion," *Science*, vol. 178, pp. 178–179, Oct. 1972.

- [347] J. Cha, M. Eid, L. Rahal, and A. E. Saddik, "HugMe: an interpersonal haptic communication system," in *Haptic Audio visual Environments and Games, 2008. HAVE 2008. IEEE International Workshop on*, pp. 99–102, 2008.
- [348] J. B. F. V. Erp, H. A. H. C. V. Veen, C. Jansen, and T. Dobbins, "Waypoint navigation with a vibrotactile waist belt," *ACM Trans. Appl. Percept.*, vol. 2, no. 2, pp. 106–117, 2005.
- [349] W. Heuten, N. Henze, S. Boll, and M. Pielot, "Tactile wayfinder: a non-visual support system for wayfinding," in *Proceedings of the 5th Nordic conference on Human-computer interaction: building bridges*, (Lund, Sweden), pp. 172–181, ACM, 2008.
- [350] K. Tsukada and M. Yasumura, "ActiveBelt: Belt-Type wearable tactile display for directional navigation," in *UbiComp 2004: Ubiquitous Computing*, pp. 399, 384, 2004.
- [351] C. Wall and M. Weinberg, "Balance prostheses for postural control," *Engineering in Medicine and Biology Magazine, IEEE*, vol. 22, no. 2, pp. 84–90, 2003.
- [352] R. W. Lindeman, Y. Yanagida, H. Noma, and K. Hosaka, "Wearable vibrotactile systems for virtual contact and information display," *Virtual Reality*, vol. 9, no. 2-3, pp. 203–213, 2005.
- [353] A. H. Rupert, "An instrumentation solution for reducing spatial disorientation mishaps," *IEEE Engineering in Medicine and Biology Magazine: The Quarterly Magazine of the Engineering in Medicine & Biology Society*, vol. 19, pp. 71–80, Apr. 2000. PMID: 10738664.
- [354] J. van Erp and H. A. van Veen, "A multi-purpose tactile vest for astronauts in the international space station," pp. 405–408, In *Proceedings of Eurohaptics 2003*, 2003.
- [355] R. Cholewiak, C. Brill, and A. Schwab, "Vibrotactile localization on the abdomen: Effects of place and space," *Perception & Psychophysics*, vol. 66, no. 6, pp. 987, 970, 2004.
- [356] A. Ferscha, B. Emsenhuber, A. Riener, C. Holzman, M. Hechinger, D. Hochreiter, M. Franz, D. Zeider, M. D. S. Rocha, and C. Klein, "Video paper: Vibro-tactile space awareness," *Adjunct Proceedings of the 10th International Conference on Ubiquitous Computing*, 2008.
- [357] S. Ram and J. Sharf, "The people sensor: a mobility aid for the visually impaired," pp. 166–167, 1998.
- [358] I. H. Computer, "Guidelines for the use of Vibro-Tactile displays," 2002.

- [359] J. Ryu and S. Choi, “posvibeditor: Graphical authoring tool of vibrotactile patterns,” in *Haptic Audio visual Environments and Games, 2008. HAVE 2008. IEEE International Workshop on*, pp. 120–125, 2008.
- [360] P. Barralon, G. Ng, G. Dumont, S. K. W. Schwarz, and M. Ansermino, “Development and evaluation of multidimensional tactons for a wearable tactile display,” in *Proceedings of the 9th international conference on Human computer interaction with mobile devices and services*, (Singapore), pp. 186–189, ACM, 2007.
- [361] L. Brown, S. Brewster, and H. Purchase, “A first investigation into the effectiveness of tactons,” in *Eurohaptics Conference, 2005 and Symposium on Haptic Interfaces for Virtual Environment and Teleoperator Systems, 2005. World Haptics 2005. First Joint*, pp. 167–176, 2005.
- [362] H. A. H. C. v. Veen and J. B. F. v. Erp, “Tactile information presentation in the cockpit,” in *Proceedings of the First International Workshop on Haptic Human-Computer Interaction*, (London, UK), pp. 174–181, Springer-Verlag, 2001.
- [363] R. W. Picard, *Affective Computing*. MIT Press, 2000.
- [364] T. L. McDaniel, S. Krishna, D. Colbry, and S. Panchanathan, “Using tactile rhythm to convey interpersonal distances to individuals who are blind,” in *Proceedings of the 27th international conference extended abstracts on Human factors in computing systems*, (Boston, MA, USA), pp. 4669–4674, ACM, 2009.
- [365] N. Edwards, J. Rosenthal, D. Molbery, J. Lindsey, K. Blair, T. McDaniel, S. Krishna, and S. Panchanathan, “A pragmatic approach to the design and implementation of a vibrotactile belt and its applications,” (Italy), pp. 13–18, 2009.
- [366] J. Lieberman and C. Breazeal, “Development of a wearable vibrotactile feedback suit for accelerated human motor learning,” pp. 4001–4006, 2007.
- [367] R. Lindeman, Y. Yanagida, K. Hosaka, and S. Abe, “The TactaPack: a wireless Sensor/Actuator package for physical therapy applications,” in *Haptic Interfaces for Virtual Environment and Teleoperator Systems, 2006 14th Symposium on*, pp. 337–341, 2006.
- [368] D. Drobny, M. Weiss, and J. Borchers, “Saltate!: a sensor-based system to support dance beginners,” in *Proceedings of the 27th international conference extended abstracts on Human factors in computing systems - CHI EA '09*, (Boston, MA, USA), p. 3943, 2009.

- [369] D. Spelmezan, M. Jacobs, A. Hilgers, and J. Borchers, “Tactile motion instructions for physical activities,” in *Proceedings of the 27th international conference on Human factors in computing systems - CHI '09*, (Boston, MA, USA), p. 2243, 2009.
- [370] J. V. B. Linden, E. Schoonderwaldt, and J. Bird, “Good vibrations: guiding body movements with vibrotactile feedback,” (Cambridge), 2009.
- [371] T. Grosshauser and T. Hermann, “Augmented haptics an interactive feedback system for musicians,” *Haptic and Audio Interaction Design*, vol. 5763/2009, pp. 100–108, 2009.
- [372] K. Frster, M. Bchlin, and G. Trster, “Non-interrupting user interfaces for electronic body-worn swim devices,” (Corfu, Greece), pp. 1–4, ACM, 2009.
- [373] A. Nakamura, S. Tabata, T. Ueda, S. Kiyofuji, and Y. Kuno, “Multimodal presentation method for a dance training system,” in *CHI '05 extended abstracts on Human factors in computing systems - CHI '05*, (Portland, OR, USA), p. 1685, 2005.
- [374] H. E. Gardner, *Frames Of Mind: The Theory Of Multiple Intelligences*. Basic Books, 10th ed., Apr. 1993.
- [375] T. Bradberry and J. Greaves, *Emotional Intelligence 2.0*. TalentSmart, Har/Onl en ed., June 2009.
- [376] E. L. Thorndike, “Intelligence and its uses,” *Harper’s Magazine*, vol. 140, p. 227235, 1920.
- [377] K. Albrecht, *Social Intelligence: The New Science of Success*. Pfeiffer, Nov. 2005.
- [378] G. Matthews, M. Zeidner, and R. D. Roberts, *Science of Emotional Intelligence: Knowns and Unknowns*. Oxford University Press, USA, 1 ed., Aug. 2007.
- [379] D. Goleman, *Working with Emotional Intelligence*. Bantam, Jan. 2000.
- [380] K. V. Petrides, R. Pita, and F. Kokkinaki, “The location of trait emotional intelligence in personality factor space,” *British Journal of Psychology*, vol. 98, pp. 273–289, May 2007.
- [381] N. K. Humphrey, “Vision in a monkey without striate cortex: a case study,” *Perception*, vol. 3, pp. 241–255, 1974.
- [382] L. Brothers, “The social brain: A project for integrating primate behavior and neurophysiology in a new domain.,” *Concepts in Neuroscience*, vol. 1, pp. 51, 27.

- [383] A. Beck, C. Ward, M. Mendelson, J. Mock, and J. Erbaugh, "An inventory for measuring depression," *Archives of General Psychiatry*, vol. 4, pp. 571, 561, June 1961.
- [384] D. W. Russell, "UCLA loneliness scale (Version 3): reliability, validity, and factor structure," *Journal of Personality Assessment*, vol. 66, pp. 20–40, Feb. 1996. PMID: 8576833.
- [385] R. E. Riggio, *Social Skills Inventory*. Palo Alto, CA: Consulting Psychologists Press, 1989.
- [386] R. E. Riggio and J. Zimmermann, "Social skills and interpersonal competence: Influences on social support and social seeking," in *Advances in Personal Relationships* (W. H. Jones and D. Perlman, eds.), pp. 133–155, London: Jessica Kingsley, 1991.
- [387] D. Magnusson, "An analysis of situational dimensions," *Perceptual and Motor Skills*, vol. 32, pp. 851–867, 1991.
- [388] W. Barfield, "Information privacy as a function of facial recognition technology and wearable computers," *bepress Legal Repository*, vol. 1739, pp. 1–75, 2006.
- [389] M. G. Morgan and E. Newton, "Protecting public anonymity," *Issues in Science and Technology*, 2004.

APPENDIX A

ALGORITHM FOR ESTIMATING RANK AVERAGE OF GROUPS

While analyzing the responses of participants to the online survey, the participants responses for each question are represented as entries $x_{i,q}$, where, i represents the i^{th} participant and q represents the q^{th} question. $i = 1, \dots, N$ are the N participants who responded on the survey, and $q = 1, \dots, Q$ are the Q questions. In the survey presented in Chapter 3, $N = 28$ and $Q = 8$.

A.0.1 Procedure

Input: Each participants response is considered as an entry e_m into a pool $E = \{x_{i,q}\}$, where, $m = 1, \dots, M$, and $M = N \times Q$.

Output: The rank average for the Q groups (questions), \bar{R}_m .

Steps:

1. Group $e_n \in E$ removing all group affiliations.
2. Order the entries from 1 to M and assign a rank r_{iq} .
3. Assign any tied values the average of the ranks they would have received had they not been tied.
4. Rank Average for each group is then given as

$$\bar{R}_m = \frac{\sum_{i \in Q_m, q=m} r_{iq}}{n_m} \quad (\text{A.1})$$

Where, Q_m represents the group m with the cardinality n_m .

Since no assumptions on the distribution of the response are made, unlike the mean, the rank average gives a non-parametric method for comparing the groups.

APPENDIX B

CONVEX OPTIMIZATION USING NEWTON'S METHOD - ENTROPY MAXIMIZATION UNDER LINEAR CONSTRAINTS

B.1 Entropy Maximization under Linear Constraints Problem:

$$\text{Maximize } f(x_i) = \left(- \sum_{i=1}^N x_i \ln(x_i) \right)$$

$$\text{Subject to } g_1(x_i) = \sum_{i=1}^N x_i = 1$$

$$g_2(x_i) = \sum_{i=1}^N x_i l_i = L$$

Given the objective and the constraints, Lagrange Multipliers is a popular method for combining the objective and constraints into a single optimization strategy with a set of unknowns. If $X = \{x_i\}$ represents the optimization vector, Lagrange Multipliers rely on the fact that at the optimal solution $X^* = \{x_i^*\}$, the gradients of the objective function and the constraint should be proportional. Figure B.1 illustrates a problem of minimizing objective $f(x,y)$ under the constraint of $g(x,y) = c$.

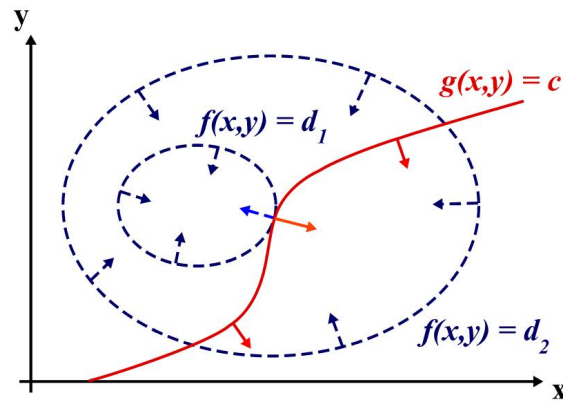


Figure B.1: Find x and y to maximize $f(x,y)$ subject to a constraint $g(x,y) = c$

The Lagrange Function $\Lambda(X, \alpha, \beta)$ becomes,

$$\Lambda(X, \alpha, \beta) = \left(- \sum_{i=1}^N x_i \ln(x_i) \right) + \alpha \left(\sum_{i=1}^N x_i - 1 \right) + \beta \left(\sum_{i=1}^N x_i l_i - L \right) \quad (\text{B.1})$$

Finding the various gradients, we have

$$\frac{\delta f(x_i)}{\delta x_i} = \ln\left(\frac{1}{x_i}\right) + 1 \quad (\text{B.2})$$

$$\frac{\delta g_1(x_i)}{\delta x_i} = 1 \quad (\text{B.3})$$

$$\frac{\delta g_2(x_i)}{\delta x_i} = l_i \quad (\text{B.4})$$

From the theory of Lagrange Multipliers, we know that the gradient of the objective is equal to the gradients of the constraints at the optimal solution x_i^* . Thus,

$$\nabla f(x_i^*) = \alpha \nabla g_1(x_i^*) + \beta \nabla g_2(x_i^*) \quad (\text{B.5})$$

$$g_1(x_i^*) = 1 \quad (\text{B.6})$$

$$g_2(x_i^*) = L \quad (\text{B.7})$$

which can be rewritten as,

$$\ln\left(\frac{1}{x_i^*} + 1\right) = \alpha + \beta l_i, i = \{1, \dots, N\} \quad (\text{B.8})$$

$$\sum_{i=1}^N x_i^* = 1 \quad (\text{B.9})$$

$$\sum_{i=1}^N x_i^* l_i = L \quad (\text{B.10})$$

From B.8 we can write

$$x_i^* = e^{1-\alpha-\beta l_i} = e^{1-\alpha} e^{-\beta l_i} \quad (\text{B.11})$$

Equation B.11 can be rewritten as

$$x_i^* = \frac{e^{1-\alpha} e^{-\beta l_i}}{1} \quad (\text{B.12})$$

$$x_i^* = \frac{e^{1-\alpha} e^{-\beta l_i}}{\sum_{i=1}^N x_i^*} \quad \text{from Equation B.9} \quad (\text{B.13})$$

$$x_i^* = \frac{e^{1-\alpha} e^{-\beta l_i}}{\sum_{i=1}^N e^{1-\alpha} e^{-\beta l_i}} \quad (\text{B.14})$$

$$x_i^* = \frac{e^{-\beta l_i}}{\sum_{i=1}^N e^{-\beta l_i}} \quad (\text{B.15})$$

This eliminates one set of Lagrange Multipliers, α . Now from Equation B.10

$$\sum_{i=1}^N x_i^* l_i = L \quad (\text{B.16})$$

$$\sum_{i=1}^N \left(\frac{e^{-\beta l_i}}{\sum_{j=1}^N e^{-\beta l_j}} \right) l_i = L \quad (\text{B.17})$$

$$\frac{\sum_{i=1}^N e^{-\beta l_i} l_i}{\sum_{j=1}^N e^{-\beta l_j}} = L \quad (\text{B.18})$$

$$\sum_{i=1}^N e^{-\beta l_i} = \sum_{j=1}^N e^{-\beta l_j} L \quad (\text{B.19})$$

$$\sum_{i=1}^N e^{-\beta l_i} (l_i - L) = 0 \quad (\text{B.20})$$

From Equation B.20 we identify the Lagrange Multipliers to be the β s and the solution to the required x_i s can then be found through Equation B.15.

The β s can be found through numerical optimization and in the section below, we show the use of Newton's method for accomplishing the same.

B.2 Newton's method

Newton's method is a popular numerical optimization technique used to find the roots of an function $f(\mathbf{x}) = 0$.

From Equation B.20, we know that the β s form the root of the equation

$$f(\beta) = \sum_{i=1}^N e^{-\beta l_i} (l_i - L) = 0 \quad (\text{B.21})$$

At each iteration the method approximates $f(\beta)$ by a quadratic function around β , and then takes a step towards the maximum/minimum of that quadratic function. Typically, the quadratic approximation is obtained from the function itself through its second order Taylor expansion.

$$f(\beta + \nabla\beta) = f(\beta) + f'(\beta)\Delta\beta + \frac{1}{2}f''(\beta)(\Delta\beta)^2$$

Newton's method determines the next approximation for the roots, β_{n+1} , by constructing a tangent to the objective function, $f(\beta)$, at the point β_n and determining the tangent's roots, i.e. determining the intersection of the tangent, $f'(\beta_n)$, with the β axis.

Thus, the next iteration for β_n could be obtained by determining $\Delta\beta$ by setting the derivative of the objective function to 0.

$$f'(\beta_n) = 0 \tag{B.22}$$

From the Taylor Expansion,

$$\begin{aligned} f'(\beta_n + \Delta\beta) &= f'(\beta_n) + \Delta\beta f''(\beta_n) \\ f'(\beta_n) + \Delta\beta f''(\beta_n) &= 0 \\ \Delta\beta &= -\frac{f'(\beta_n)}{f''(\beta_n)} \end{aligned}$$

Thus, the next best approximation for the roots can be obtained as,

$$\begin{aligned} \beta_{n+1} &= \beta_n + \Delta\beta \\ \beta_{n+1} &= \beta_n - \frac{f'(\beta_n)}{f''(\beta_n)} \end{aligned}$$

The iterations are continued until the optimal solution is found within a numerical error bound or until a predetermined number of iterations have been completed.

B.2.1 Newton's Method for Estimating Weights w_{ij}

Equation B.20 has a similar form as the Equation 5.11. The steps below shows the Newton method optimization to determine the w_{ij} s.

Step 0: Initialize

$$\left. \begin{aligned} \Phi &= \begin{pmatrix} 1 & \dots & 1 & \dots & 1 & \dots & 1 \\ L_{11} & \dots & L_{1k} & \dots & L_{N1} & \dots & L_{Nk} \end{pmatrix} \\ \beta_{N \times k+1} &= [0, \dots, 0]^T \\ x &= [1 \quad X(t)]^T \\ \varepsilon &= 10^{-3} \end{aligned} \right\} \quad (\text{B.23})$$

Step 1

$$\left. \begin{aligned} u_i &= \exp \beta^T \Phi(:, i) \\ \text{where,} \\ u &= [u_1, u_2, \dots, u_{N \times k}]^T \end{aligned} \right\} \quad (\text{B.24})$$

Step 2

$$\left. \begin{aligned} v_j &= x_j - \sum_{m=1}^{N \times k} \Phi(j, m) u_m \\ v &= [v_1, v_2, \dots, v_{D+1}] \\ \text{where, } D &\text{ is the dimension of the data} \end{aligned} \right\} \quad (\text{B.25})$$

Step 3

$$\left. \begin{aligned} g_{ij} &= \sum_{m=1}^{N \times k} \Phi(i, m) \Phi(j, m) u_m \\ \text{where, } G &\text{ is the Jacobian Matrix} \end{aligned} \right\} \quad (\text{B.26})$$

Step 4

$$\beta_{New} = \beta_{Old} + G^{-1}v \quad (\text{B.27})$$

Step 5

$$\text{If } \|G^{-1}v\| > \varepsilon, \text{ go back to Step 1.} \quad (\text{B.28})$$

APPENDIX C

AMERICAN COMMUNITY SURVEY FORM - SAMPLE PAGES FROM 2008

SURVEY FORM



U.S. DEPARTMENT OF COMMERCE
Economics and Statistics Administration
U.S. CENSUS BUREAU

THE American Community Survey

This booklet shows the content of the American Community Survey questionnaire.

Please complete this form and return it as soon as possible after receiving it in the mail.

This form asks for information about the people who are living or staying at the address on the mailing label and about the house, apartment, or mobile home located at the address on the mailing label.



If you need help or have questions about completing this form, please call **1-800-354-7271**. The telephone call is free.

Telephone Device for the Deaf (TDD):
Call 1-800-582-8330. The telephone call is free.

¿NECESITA AYUDA? Si usted habla español y necesita ayuda para completar su cuestionario, llame sin cargo alguno al **1-877-833-5625**.

Usted también puede pedir un cuestionario en español o completar su entrevista por teléfono con un entrevistador que habla español.

For more information about the American Community Survey, visit our web site at: <http://www.census.gov/acs/www/>

U S C E N S U S B U R E A U



Start Here

➔ **Please print today's date.**

Month Day Year

➔ **Please print the name and telephone number of the person who is filling out this form.** We may contact you if there is a question.

Last Name

First Name MI

Area Code + Number
 -

➔ **How many people are living or staying at this address?**

- **INCLUDE** everyone who is living or staying here for more than 2 months.
- **INCLUDE** yourself if you are living here for more than 2 months.
- **INCLUDE** anyone else staying here who does not have another place to stay, even if they are here for 2 months or less.
- **DO NOT INCLUDE** anyone who is living somewhere else for more than 2 months, such as a college student living away or someone in the Armed Forces on deployment.

Number of people

➔ **Fill out pages 2, 3, and 4 for everyone, including yourself, who is living or staying at this address for more than 2 months. Then complete the rest of the form.**

FORM **ACS-1 (INFO) (2010) KFI**
(09-14-2009)

OMB No. 0607-0810

Person 1	Person 2
<p>(Person 1 is the person living or staying here in whose name this house or apartment is owned, being bought, or rented. If there is no such person, start with the name of any adult living or staying here.)</p>	
<p>1 What is Person 1's name? Last Name (Please print) <input type="text"/> First Name <input type="text"/> MI <input type="text"/></p>	<p>1 What is Person 2's name? Last Name (Please print) <input type="text"/> First Name <input type="text"/> MI <input type="text"/></p>
<p>2 How is this person related to Person 1? <input checked="" type="checkbox"/> Person 1 <input type="checkbox"/> Husband or wife <input type="checkbox"/> Biological son or daughter <input type="checkbox"/> Adopted son or daughter <input type="checkbox"/> Stepson or stepdaughter <input type="checkbox"/> Brother or sister <input type="checkbox"/> Father or mother <input type="checkbox"/> Grandchild <input type="checkbox"/> Parent-in-law</p>	<p>2 How is this person related to Person 1? Mark (X) ONE box. <input type="checkbox"/> Son-in-law or daughter-in-law <input type="checkbox"/> Other relative <input type="checkbox"/> Roomer or boarder <input type="checkbox"/> Housemate or roommate <input type="checkbox"/> Unmarried partner <input type="checkbox"/> Foster child <input type="checkbox"/> Other nonrelative</p>
<p>3 What is Person 1's sex? Mark (X) ONE box. <input type="checkbox"/> Male <input type="checkbox"/> Female</p>	<p>3 What is Person 2's sex? Mark (X) ONE box. <input type="checkbox"/> Male <input type="checkbox"/> Female</p>
<p>4 What is Person 1's age and what is Person 1's date of birth? Please report babies as age 0 when the child is less than 1 year old. Print numbers in boxes. Age (in years) <input type="text"/> <input type="text"/> <input type="text"/> <input type="text"/> Month <input type="text"/> Day <input type="text"/> Year of birth <input type="text"/> <input type="text"/> <input type="text"/> <input type="text"/></p>	<p>4 What is Person 2's age and what is Person 2's date of birth? Please report babies as age 0 when the child is less than 1 year old. Print numbers in boxes. Age (in years) <input type="text"/> <input type="text"/> <input type="text"/> <input type="text"/> Month <input type="text"/> Day <input type="text"/> Year of birth <input type="text"/> <input type="text"/> <input type="text"/> <input type="text"/></p>
<p>→ NOTE: Please answer BOTH Question 5 about Hispanic origin and Question 6 about race. For this survey, Hispanic origins are not races.</p>	
<p>5 Is Person 1 of Hispanic, Latino, or Spanish origin? <input type="checkbox"/> No, not of Hispanic, Latino, or Spanish origin <input type="checkbox"/> Yes, Mexican, Mexican Am., Chicano <input type="checkbox"/> Yes, Puerto Rican <input type="checkbox"/> Yes, Cuban <input type="checkbox"/> Yes, another Hispanic, Latino, or Spanish origin – Print origin, for example, Argentinean, Colombian, Dominican, Nicaraguan, Salvadoran, Spaniard, and so on. ↴</p>	<p>5 Is Person 2 of Hispanic, Latino, or Spanish origin? <input type="checkbox"/> No, not of Hispanic, Latino, or Spanish origin <input type="checkbox"/> Yes, Mexican, Mexican Am., Chicano <input type="checkbox"/> Yes, Puerto Rican <input type="checkbox"/> Yes, Cuban <input type="checkbox"/> Yes, another Hispanic, Latino, or Spanish origin – Print origin, for example, Argentinean, Colombian, Dominican, Nicaraguan, Salvadoran, Spaniard, and so on. ↴</p>
<p>6 What is Person 1's race? Mark (X) one or more boxes. <input type="checkbox"/> White <input type="checkbox"/> Black, African Am., or Negro <input type="checkbox"/> American Indian or Alaska Native – Print name of enrolled or principal tribe. ↴ <input type="checkbox"/> Asian Indian <input type="checkbox"/> Chinese <input type="checkbox"/> Filipino <input type="checkbox"/> Other Asian – Print race, for example, Hmong, Laotian, Thai, Pakistani, Cambodian, and so on. ↴ <input type="checkbox"/> Japanese <input type="checkbox"/> Korean <input type="checkbox"/> Vietnamese <input type="checkbox"/> Native Hawaiian <input type="checkbox"/> Guamanian or Chamorro <input type="checkbox"/> Samoan <input type="checkbox"/> Other Pacific Islander – Print race, for example, Fijian, Tongan, and so on. ↴ <input type="checkbox"/> Some other race – Print race. ↴</p>	<p>6 What is Person 2's race? Mark (X) one or more boxes. <input type="checkbox"/> White <input type="checkbox"/> Black, African Am., or Negro <input type="checkbox"/> American Indian or Alaska Native – Print name of enrolled or principal tribe. ↴ <input type="checkbox"/> Asian Indian <input type="checkbox"/> Chinese <input type="checkbox"/> Filipino <input type="checkbox"/> Other Asian – Print race, for example, Hmong, Laotian, Thai, Pakistani, Cambodian, and so on. ↴ <input type="checkbox"/> Japanese <input type="checkbox"/> Korean <input type="checkbox"/> Vietnamese <input type="checkbox"/> Native Hawaiian <input type="checkbox"/> Guamanian or Chamorro <input type="checkbox"/> Samoan <input type="checkbox"/> Other Pacific Islander – Print race, for example, Fijian, Tongan, and so on. ↴ <input type="checkbox"/> Some other race – Print race. ↴</p>



Housing

→ Please answer the following questions about the house, apartment, or mobile home at the address on the mailing label.

- 1 Which best describes this building?**
Include all apartments, flats, etc., even if vacant.
- A mobile home
 - A one-family house detached from any other house
 - A one-family house attached to one or more houses
 - A building with 2 apartments
 - A building with 3 or 4 apartments
 - A building with 5 to 9 apartments
 - A building with 10 to 19 apartments
 - A building with 20 to 49 apartments
 - A building with 50 or more apartments
 - Boat, RV, van, etc.

- 2 About when was this building first built?**
- 2000 or later – *Specify year* →
- 1990 to 1999
 - 1980 to 1989
 - 1970 to 1979
 - 1960 to 1969
 - 1950 to 1959
 - 1940 to 1949
 - 1939 or earlier

- 3 When did PERSON 1 (listed on page 2) move into this house, apartment, or mobile home?**
- Month Year

A Answer questions 4 – 6 if this is a HOUSE OR A MOBILE HOME; otherwise, SKIP to question 7a.

- 4 How many acres is this house or mobile home on?**
- Less than 1 acre → SKIP to question 6
 - 1 to 9.9 acres
 - 10 or more acres

- 5 IN THE PAST 12 MONTHS, what were the actual sales of all agricultural products from this property?**
- None
 - \$1 to \$999
 - \$1,000 to \$2,499
 - \$2,500 to \$4,999
 - \$5,000 to \$9,999
 - \$10,000 or more

- 6 Is there a business (such as a store or barber shop) or a medical office on this property?**
- Yes
 - No

- 7 a. How many separate rooms are in this house, apartment, or mobile home?**
Rooms must be separated by built-in archways or walls that extend out at least 6 inches and go from floor to ceiling.
- INCLUDE bedrooms, kitchens, etc.
 - EXCLUDE bathrooms, porches, balconies, foyers, halls, or unfinished basements.

Number of rooms

- b. How many of these rooms are bedrooms?**
Count as bedrooms those rooms you would list if this house, apartment, or mobile home were for sale or rent. If this is an efficiency/studio apartment, print "0".

Number of bedrooms

- 8 Does this house, apartment, or mobile home have –**
- | | Yes | No |
|--|--------------------------|--------------------------|
| a. hot and cold running water? | <input type="checkbox"/> | <input type="checkbox"/> |
| b. a flush toilet? | <input type="checkbox"/> | <input type="checkbox"/> |
| c. a bathtub or shower? | <input type="checkbox"/> | <input type="checkbox"/> |
| d. a sink with a faucet? | <input type="checkbox"/> | <input type="checkbox"/> |
| e. a stove or range? | <input type="checkbox"/> | <input type="checkbox"/> |
| f. a refrigerator? | <input type="checkbox"/> | <input type="checkbox"/> |
| g. telephone service from which you can both make and receive calls? <i>Include cell phones.</i> | <input type="checkbox"/> | <input type="checkbox"/> |

- 9 How many automobiles, vans, and trucks of one-ton capacity or less are kept at home for use by members of this household?**
- None
 - 1
 - 2
 - 3
 - 4
 - 5
 - 6 or more

- 10 Which FUEL is used MOST for heating this house, apartment, or mobile home?**
- Gas: from underground pipes serving the neighborhood
 - Gas: bottled, tank, or LP
 - Electricity
 - Fuel oil, kerosene, etc.
 - Coal or coke
 - Wood
 - Solar energy
 - Other fuel
 - No fuel used



Person 1
13190087

7 Please copy the name of Person 1 from page 2, then continue answering questions below.

Last Name

First Name MI

7 Where was this person born?

In the United States – Print name of state.

Outside the United States – Print name of foreign country, or Puerto Rico, Guam, etc.

8 Is this person a citizen of the United States?

Yes, born in the United States → SKIP to 10a

Yes, born in Puerto Rico, Guam, the U.S. Virgin Islands, or Northern Marianas

Yes, born abroad of U.S. citizen parent or parents

Yes, U.S. citizen by naturalization – Print year of naturalization

No, not a U.S. citizen

9 When did this person come to live in the United States? Print numbers in boxes.

Year

10 a. At any time IN THE LAST 3 MONTHS, has this person attended school or college? Include only nursery or preschool, kindergarten, elementary school, home school, and schooling which leads to a high school diploma or a college degree.

No, has not attended in the last 3 months → SKIP to question 11

Yes, public school, public college

Yes, private school, private college, home school

b. What grade or level was this person attending? Mark (X) ONE box.

Nursery school, preschool

Kindergarten

Grade 1 through 12 – Specify grade 1 – 12

College undergraduate years (freshman to senior)

Graduate or professional school beyond a bachelor's degree (for example: MA or PhD program, or medical or law school)

11 What is the highest degree or level of school this person has COMPLETED? Mark (X) ONE box. If currently enrolled, mark the previous grade or highest degree received.

NO SCHOOLING COMPLETED

No schooling completed

NURSERY OR PRESCHOOL THROUGH GRADE 12

Nursery school

Kindergarten

Grade 1 through 11 – Specify grade 1 – 11

12th grade – **NO DIPLOMA**

HIGH SCHOOL GRADUATE

Regular high school diploma

GED or alternative credential

COLLEGE OR SOME COLLEGE

Some college credit, but less than 1 year of college credit

1 or more years of college credit, no degree

Associate's degree (for example: AA, AS)

Bachelor's degree (for example: BA, BS)

AFTER BACHELOR'S DEGREE

Master's degree (for example: MA, MS, MEng, MEd, MSW, MBA)

Professional degree beyond a bachelor's degree (for example: MD, DDS, DVM, LLB, JD)

Doctorate degree (for example: PhD, EdD)

F Answer question 12 if this person has a bachelor's degree or higher. Otherwise, SKIP to question 13.

12 This question focuses on this person's BACHELOR'S DEGREE. Please print below the specific major(s) of any BACHELOR'S DEGREES this person has received. (For example: chemical engineering, elementary teacher education, organizational psychology)

13 What is this person's ancestry or ethnic origin?

(For example: Italian, Jamaican, African Am., Cambodian, Cape Verdean, Norwegian, Dominican, French Canadian, Haitian, Korean, Lebanese, Polish, Nigerian, Mexican, Taiwanese, Ukrainian, and so on.)

14 a. Does this person speak a language other than English at home?

Yes

No → SKIP to question 15a

b. What is this language?

(For example: Korean, Italian, Spanish, Vietnamese)

c. How well does this person speak English?

Very well

Well

Not well

Not at all

15 a. Did this person live in this house or apartment 1 year ago?

Person is under 1 year old → SKIP to question 16

Yes, this house → SKIP to question 16

No, outside the United States and Puerto Rico – Print name of foreign country, or U.S. Virgin Islands, Guam, etc., below; then SKIP to question 16

No, different house in the United States or Puerto Rico

b. Where did this person live 1 year ago?

Address (Number and street name)

Name of city, town, or post office

Name of U.S. county or municipio in Puerto Rico

Name of U.S. state or Puerto Rico ZIP Code



Person 1 (continued)

16 Is this person CURRENTLY covered by any of the following types of health insurance or health coverage plans? Mark "Yes" or "No" for EACH type of coverage in items a – h.

- | | | |
|---|-----|----|
| a. Insurance through a current or former employer or union (of this person or another family member) | Yes | No |
| b. Insurance purchased directly from an insurance company (by this person or another family member) | | |
| c. Medicare, for people 65 and older, or people with certain disabilities | | |
| d. Medicaid, Medical Assistance, or any kind of government-assistance plan for those with low incomes or a disability | | |
| e. TRICARE or other military health care | | |
| f. VA (including those who have ever used or enrolled for VA health care) | | |
| g. Indian Health Service | | |
| h. Any other type of health insurance or health coverage plan – Specify → | | |

17 a. Is this person deaf or does he/she have serious difficulty hearing?

- Yes
 No

b. Is this person blind or does he/she have serious difficulty seeing even when wearing glasses?

- Yes
 No

G Answer question 18a – c if this person is 5 years old or over. Otherwise, SKIP to the questions for Person 2 on page 12.

18 a. Because of a physical, mental, or emotional condition, does this person have serious difficulty concentrating, remembering, or making decisions?

- Yes
 No

b. Does this person have serious difficulty walking or climbing stairs?

- Yes
 No

c. Does this person have difficulty dressing or bathing?

- Yes
 No

H Answer question 19 if this person is 15 years old or over. Otherwise, SKIP to the questions for Person 2 on page 12.

19 Because of a physical, mental, or emotional condition, does this person have difficulty doing errands alone such as visiting a doctor's office or shopping?

- Yes
 No

20 What is this person's marital status?

- Now married
 Widowed
 Divorced
 Separated
 Never married → SKIP to **I**

21 In the PAST 12 MONTHS did this person get:

- | | | |
|--------------|--------------------------|--------------------------|
| | Yes | No |
| a. Married? | <input type="checkbox"/> | <input type="checkbox"/> |
| b. Widowed? | <input type="checkbox"/> | <input type="checkbox"/> |
| c. Divorced? | <input type="checkbox"/> | <input type="checkbox"/> |

22 How many times has this person been married?

- Once
 Two times
 Three or more times

23 In what year did this person last get married?

Year

I Answer question 24 if this person is female and 15 – 50 years old. Otherwise, SKIP to question 25a.

24 Has this person given birth to any children in the past 12 months?

- Yes
 No

25 a. Does this person have any of his/her own grandchildren under the age of 18 living in this house or apartment?

- Yes
 No → SKIP to question 26

b. Is this grandparent currently responsible for most of the basic needs of any grandchild(ren) under the age of 18 who live(s) in this house or apartment?

- Yes
 No → SKIP to question 26

c. How long has this grandparent been responsible for the(ese) grandchild(ren)?

If the grandparent is financially responsible for more than one grandchild, answer the question for the grandchild for whom the grandparent has been responsible for the longest period of time.

- Less than 6 months
 6 to 11 months
 1 or 2 years
 3 or 4 years
 5 or more years

26 Has this person ever served on active duty in the U.S. Armed Forces, military Reserves, or National Guard? Active duty does not include training for the Reserves or National Guard, but DOES include activation, for example, for the Persian Gulf War.

- Yes, now on active duty
 Yes, on active duty during the last 12 months, but not now
 Yes, on active duty in the past, but not during the last 12 months
 No, training for Reserves or National Guard only → SKIP to question 28a
 No, never served in the military → SKIP to question 28a

27 When did this person serve on active duty in the U.S. Armed Forces? Mark (X) a box for EACH period in which this person served, even if just for part of the period.

- September 2001 or later
 August 1990 to August 2001 (including Persian Gulf War)
 September 1980 to July 1990
 May 1975 to August 1980
 Vietnam era (August 1964 to April 1975)
 March 1961 to July 1964
 February 1955 to February 1961
 Korean War (July 1950 to January 1955)
 January 1947 to June 1950
 World War II (December 1941 to December 1946)
 November 1941 or earlier

28 a. Does this person have a VA service-connected disability rating?

- Yes (such as 0%, 10%, 20%, ... , 100%)
 No → SKIP to question 29a

b. What is this person's service-connected disability rating?

- 0 percent
 10 or 20 percent
 30 or 40 percent
 50 or 60 percent
 70 percent or higher



Person 1 (continued)

29 a. **LAST WEEK, did this person work for pay at a job (or business)?**

Yes → SKIP to question 30

No – Did not work (or retired)

b. **LAST WEEK, did this person do ANY work for pay, even for as little as one hour?**

Yes

No → SKIP to question 35a

30 **At what location did this person work LAST WEEK?** *If this person worked at more than one location, print where he or she worked most last week.*

a. **Address (Number and street name)**

If the exact address is not known, give a description of the location such as the building name or the nearest street or intersection.

b. **Name of city, town, or post office**

c. **Is the work location inside the limits of that city or town?**

Yes

No, outside the city/town limits

d. **Name of county**

e. **Name of U.S. state or foreign country**

f. **ZIP Code**

31 **How did this person usually get to work LAST WEEK?** *If this person usually used more than one method of transportation during the trip, mark (X) the box of the one used for most of the distance.*

<input type="checkbox"/> Car, truck, or van	<input type="checkbox"/> Motorcycle
<input type="checkbox"/> Bus or trolley bus	<input type="checkbox"/> Bicycle
<input type="checkbox"/> Streetcar or trolley car	<input type="checkbox"/> Walked
<input type="checkbox"/> Subway or elevated	<input type="checkbox"/> Worked at home → SKIP to question 39a
<input type="checkbox"/> Railroad	<input type="checkbox"/> Other method
<input type="checkbox"/> Ferryboat	
<input type="checkbox"/> Taxicab	

J Answer question 32 if you marked "Car, truck, or van" in question 31. Otherwise, SKIP to question 33.

32 **How many people, including this person, usually rode to work in the car, truck, or van LAST WEEK?**

Person(s)

33 **What time did this person usually leave home to go to work LAST WEEK?**

Hour Minute

a.m.

p.m.

34 **How many minutes did it usually take this person to get from home to work LAST WEEK?**

Minutes

K Answer questions 35 – 38 if this person did NOT work last week. Otherwise, SKIP to question 39a.

35 a. **LAST WEEK, was this person on layoff from a job?**

Yes → SKIP to question 35c

No

b. **LAST WEEK, was this person TEMPORARILY absent from a job or business?**

Yes, on vacation, temporary illness, maternity leave, other family/personal reasons, bad weather, etc. → SKIP to question 38

No → SKIP to question 36

c. **Has this person been informed that he or she will be recalled to work within the next 6 months OR been given a date to return to work?**

Yes → SKIP to question 37

No

36 **During the LAST 4 WEEKS, has this person been ACTIVELY looking for work?**

Yes

No → SKIP to question 38

37 **LAST WEEK, could this person have started a job if offered one, or returned to work if recalled?**

Yes, could have gone to work

No, because of own temporary illness

No, because of all other reasons (in school, etc.)

38 **When did this person last work, even for a few days?**

Within the past 12 months

1 to 5 years ago → SKIP to L

Over 5 years ago or never worked → SKIP to question 47

39 a. **During the PAST 12 MONTHS (52 weeks), did this person work 50 or more weeks? Count paid time off as work.**

Yes → SKIP to question 40

No

b. **How many weeks DID this person work, even for a few hours, including paid vacation, paid sick leave, and military service?**

50 to 52 weeks

48 to 49 weeks

40 to 47 weeks

27 to 39 weeks

14 to 26 weeks

13 weeks or less

40 **During the PAST 12 MONTHS, in the WEEKS WORKED, how many hours did this person usually work each WEEK?**

Usual hours worked each WEEK



