

Multimodal Affective Computing

Affective Information Representation,
Modelling, and Analysis

Gyanendra K. Verma



Bentham Books

**Multimodal Affective
Computing:
Affective Information
Representation, Modelling, and
Analysis**

Authored By

Gyanendra K. Verma

*Department of Information Technology, National Institute of
Technology Raipur,
Chhattisgarh, India*

**Multimodal Affective Computing:
Affective Information Representation, Modelling, and Analysis**

Author: Gyanendra K. Verma

ISBN (Online): 978-981-5124-45-3

ISBN (Print): 978-981-5124-46-0

ISBN (Paperback): 978-981-5124-47-7

© 2023, Bentham Books imprint.

Published by Bentham Science Publishers Pte. Ltd. Singapore. All Rights Reserved.

First published in 2023.

BENTHAM SCIENCE PUBLISHERS LTD.

End User License Agreement (for non-institutional, personal use)

This is an agreement between you and Bentham Science Publishers Ltd. Please read this License Agreement carefully before using the ebook/echapter/ejournal (“**Work**”). Your use of the Work constitutes your agreement to the terms and conditions set forth in this License Agreement. If you do not agree to these terms and conditions then you should not use the Work.

Bentham Science Publishers agrees to grant you a non-exclusive, non-transferable limited license to use the Work subject to and in accordance with the following terms and conditions. This License Agreement is for non-library, personal use only. For a library / institutional / multi user license in respect of the Work, please contact: permission@benthamscience.net.

Usage Rules:

1. All rights reserved: The Work is the subject of copyright and Bentham Science Publishers either owns the Work (and the copyright in it) or is licensed to distribute the Work. You shall not copy, reproduce, modify, remove, delete, augment, add to, publish, transmit, sell, resell, create derivative works from, or in any way exploit the Work or make the Work available for others to do any of the same, in any form or by any means, in whole or in part, in each case without the prior written permission of Bentham Science Publishers, unless stated otherwise in this License Agreement.
2. You may download a copy of the Work on one occasion to one personal computer (including tablet, laptop, desktop, or other such devices). You may make one back-up copy of the Work to avoid losing it.
3. The unauthorised use or distribution of copyrighted or other proprietary content is illegal and could subject you to liability for substantial money damages. You will be liable for any damage resulting from your misuse of the Work or any violation of this License Agreement, including any infringement by you of copyrights or proprietary rights.

Disclaimer:

Bentham Science Publishers does not guarantee that the information in the Work is error-free, or warrant that it will meet your requirements or that access to the Work will be uninterrupted or error-free. The Work is provided "as is" without warranty of any kind, either express or implied or statutory, including, without limitation, implied warranties of merchantability and fitness for a particular purpose. The entire risk as to the results and performance of the Work is assumed by you. No responsibility is assumed by Bentham Science Publishers, its staff, editors and/or authors for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products instruction, advertisements or ideas contained in the Work.

Limitation of Liability:

In no event will Bentham Science Publishers, its staff, editors and/or authors, be liable for any damages, including, without limitation, special, incidental and/or consequential damages and/or damages for lost data and/or profits arising out of (whether directly or indirectly) the use or inability to use the Work. The entire liability of Bentham Science Publishers shall be limited to the amount actually paid by you for the Work.

General:

1. Any dispute or claim arising out of or in connection with this License Agreement or the Work (including non-contractual disputes or claims) will be governed by and construed in accordance with the laws of Singapore. Each party agrees that the courts of the state of Singapore shall have exclusive jurisdiction to settle any dispute or claim arising out of or in connection with this License Agreement or the Work (including non-contractual disputes or claims).
2. Your rights under this License Agreement will automatically terminate without notice and without the

need for a court order if at any point you breach any terms of this License Agreement. In no event will any delay or failure by Bentham Science Publishers in enforcing your compliance with this License Agreement constitute a waiver of any of its rights.

3. You acknowledge that you have read this License Agreement, and agree to be bound by its terms and conditions. To the extent that any other terms and conditions presented on any website of Bentham Science Publishers conflict with, or are inconsistent with, the terms and conditions set out in this License Agreement, you acknowledge that the terms and conditions set out in this License Agreement shall prevail.

Bentham Science Publishers Pte. Ltd.

80 Robinson Road #02-00

Singapore 068898

Singapore

Email: subscriptions@benthamscience.net



CONTENTS

FOREWORD	i
PREFACE	ii
CONSENT FOR PUBLICATION	iii
CONFLICT OF INTEREST	iii
ACKNOWLEDGEMENTS	iv
CHAPTER 1 AFFECTIVE COMPUTING	1
1.1. INTRODUCTION	1
1.2. WHAT IS EMOTION?	2
1.2.1. Affective Human-Computer Interaction	2
1.3. BACKGROUND	3
1.4. THE ROLE OF EMOTIONS IN DECISION MAKING	4
1.5. CHALLENGES IN AFFECTIVE COMPUTING	5
1.5.1. How Can Many Emotions Be Analyzed in a Single Framework?	5
1.5.2. How Can Complex Emotions Be Represented in a Single Framework Or Model?	6
1.5.3. Is The Chosen Theoretical Viewpoint Relevant to other Areas Of Affective Computing?	6
1.5.4. How Can Physiological Signals Be Used to Anticipate Complicated Emotions?	6
1.6. AFFECTIVE COMPUTING IN PRACTICE	6
1.6.1. Avatars or Virtual Agents	7
1.6.2. Robotics	7
1.6.3. Gaming	8
1.6.4. Education	9
1.6.5. Medical	9
1.6.6. Smart Homes and Workplace Environments	10
CONCLUSION	10
REFERENCES	10
CHAPTER 2 AFFECTIVE INFORMATION REPRESENTATION	13
2.1. INTRODUCTION	13
2.2. AFFECTIVE COMPUTING AND EMOTION	13
2.2.1. Affective Human-Computer Interaction	14
2.2.2. Human Emotion Expression and Perception	15
2.2.2.1. <i>Facial Expressions</i>	15
2.2.2.2. <i>AudioHG</i>	15
2.2.2.3. <i>Physiological Signals</i>	16
2.2.2.4. <i>Hand and Gesture Movement</i>	17
2.3. RECOGNITION OF FACIAL EMOTION	17
2.3.1. Facial Expression Fundamentals	18
2.3.2. Emotion Modeling	19
2.3.3. Representation of Facial Expression	20
2.3.4. Facial Emotion's Limitations	21
2.3.5. Techniques for Classifying Facial Expressions	21
CONCLUSION	25
REFERENCES	26
CHAPTER 3 MODELS AND THEORY OF EMOTION	30
3.1. INTRODUCTION	30
3.2. EMOTION THEORY	30
3.2.1. Categorical Approach	31

3.2.2. Evolutionary Theory of Emotion by Darwin	32
3.2.3. Cognitive Appraisal and Physiological Theory of Emotions	33
3.2.4. Dimensional Approaches to Emotions	34
CONCLUSION	37
REFERENCES	37
CHAPTER 4 AFFECTIVE INFORMATION EXTRACTION, PROCESSING AND EVALUATION	40
4.1. INTRODUCTION	40
4.2. AFFECTIVE INFORMATION EXTRACTION AND PROCESSING	40
4.2.1. Information Extraction from Audio	40
4.2.2. Information Extraction from Video	41
4.2.3. Information Extraction from Physiological Signals	41
4.3. STUDIES ON AFFECT INFORMATION PROCESSING	42
4.4. EVALUATION	43
4.4.1. Types of Errors	43
4.4.1.1. <i>False Acceptance Ratio</i>	43
4.4.1.2. <i>False Reject Ratio</i>	44
4.4.2. Threshold Criteria	44
4.4.3. Performance Criteria	44
4.4.4. Evaluation Metrics	45
4.4.4.1. <i>Mean Absolute Error (MAE)</i>	45
4.4.4.2. <i>Mean Square Error (MSE)</i>	45
4.4.5. ROC Curves	45
4.4.6. F1 Measure	46
CONCLUSION	47
REFERENCES	47
CHAPTER 5 MULTIMODAL AFFECTIVE INFORMATION FUSION	49
5.1. INTRODUCTION	49
5.2. MULTIMODAL INFORMATION FUSION	49
5.2.1. Early Fusion	50
5.2.2. Intermediate Fusion	51
5.2.3. Late Fusion	51
5.3. LEVELS OF INFORMATION FUSION	53
5.3.1. Sensor or Data-level Fusion	54
5.3.2. Feature Level Fusion	55
5.3.3. Decision-Level Fusion	55
5.4. MAJOR CHALLENGES IN INFORMATION FUSION	55
CONCLUSION	56
REFERENCES	56
CHAPTER 6 MULTIMODAL FUSION FRAMEWORK AND MULTIREOLUTION ANALYSIS	59
6.1. INTRODUCTION	59
6.2. THE BENEFITS OF MULTIMODAL FEATURES	59
6.2.1. Noise In Sensed Data	60
6.2.2. Non-Universality	60
6.2.3. Complementary Information	61
6.3. FEATURE LEVEL FUSION	61
6.4. MULTIMODAL FEATURE-LEVEL FUSION	62
6.4.1. Feature Normalization	62

6.4.2. Feature Selection	63
6.4.3. Criteria For Feature Selection	63
6.5. MULTIMODAL FUSION FRAMEWORK	65
6.5.1. Feature Extraction and Selection	65
6.5.1.1. <i>Extraction of Audio Features</i>	65
6.5.1.2. <i>Extraction of Video Features</i>	65
6.5.1.3. <i>Extraction of Peripheral Features from EEG</i>	66
6.5.2. Dimension Reduction and Feature-level Fusion	66
6.5.3. Emotion Mapping to a 3D VAD Space	67
6.6. MULTIREOLUTION ANALYSIS	70
6.6.1. Motivations for the use of Multiresolution Analysis	71
6.6.2. The Wavelet Transform	71
6.6.3. The Curvelet Transform	72
6.6.4. The Ridgelet Transform	73
CONCLUSION	73
REFERENCES	73
CHAPTER 7 EMOTION RECOGNITION FROM FACIAL EXPRESSION IN A NOISY ENVIRONMENT	75
7.1. INTRODUCTION	75
7.2. THE CHALLENGES IN FACIAL EMOTION RECOGNITION	76
7.3. NOISE AND DYNAMIC RANGE IN DIGITAL IMAGES	78
7.3.1. Characteristic Sources Of Digital Image Noise	79
7.3.1.1. <i>Sensor Read Noise</i>	79
7.3.1.2. <i>Pattern Noise</i>	79
7.3.1.3. <i>Thermal Noise</i>	79
7.3.1.4. <i>Pixel Response Non-uniformity (PRNU)</i>	79
7.3.1.5. <i>Quantization Error</i>	79
7.4. THE DATABASE	80
7.4.1. Cohn-Kanade Database	80
7.4.2. JAFFE Database	80
7.4.3. In-House Database	80
7.5. EXPERIMENTS WITH THE PROPOSED FRAMEWORK	80
7.5.1. Image Pre-Processing	82
7.5.2. Feature Extraction	82
7.5.3. Feature Matching	82
7.6. RESULTS AND DISCUSSIONS	84
7.7. RESULTS UNDER ILLUMINATION CHANGES	87
7.8. RESULTS UNDER GAUSSIAN NOISE	87
7.8.1. Comparison with other Strategies	87
CONCLUSION	94
REFERENCES	94
CHAPTER 8 SPONTANEOUS EMOTION RECOGNITION FROM AUDIO-VISUAL SIGNALS	97
8.1. INTRODUCTION	97
8.2. RECOGNITION OF SPONTANEOUS EFFECTS	98
8.3. THE DATABASE	98
8.3.1. eINTERFACE Database	99
8.3.2. RML Database	100
8.4. AUDIO-BASED EMOTION RECOGNITION SYSTEM	100
8.4.1. Experiments	101

8.4.2. System Development	101
8.4.2.1. Audio Features	101
8.5. VISUAL CUE-BASED EMOTION RECOGNITION SYSTEM	104
8.5.1. Experiments	104
8.5.2. System Development	104
8.5.2.1. Visual Feature	104
8.6. EXPERIMENTS BASED ON THE PROPOSED AUDIO-VISUAL CUES FUSION	
FRAMEWORK	107
8.6.1. Results	109
8.6.2. Comparison To Other Research	110
CONCLUSION	111
REFERENCES	111
CHAPTER 9 MULTIMODAL FUSION FRAMEWORK: EMOTION RECOGNITION FROM	
PHYSIOLOGICAL SIGNALS	115
9.1. INTRODUCTION	115
9.1.1. Electrical Brain Activity	116
9.1.2. Muscle Activity	117
9.1.3. Skin Conductivity	117
9.1.4. Skin Temperature	117
9.2. MULTIMODAL EMOTION DATABASE	117
9.2.1. DEAP Database	118
9.3. FEATURE EXTRACTION	118
9.3.1. Feature Extraction from EEG	119
9.3.2. Feature Extraction from Peripheral Signals	119
9.4. CLASSIFICATION AND RECOGNITION OF EMOTION	120
9.4.1. Support Vector Machine (SVM)	120
9.4.2. Multi-Layer Perceptron (MLP)	121
9.4.3. K-Nearest Neighbor (K-NN)	122
9.5. RESULTS AND DISCUSSION	123
9.5.1. Emotion Categorization Results Based On The Proposed Multimodal Fusion	
Architecture	123
CONCLUSION	126
REFERENCES	126
CHAPTER 10 EMOTIONS MODELLING IN 3D SPACE	128
10.1. INTRODUCTION	128
10.2. AFFECT REPRESENTATION IN 2D SPACE	129
10.3. EMOTION REPRESENTATION IN 3D SPACE	131
10.4. 3D EMOTION MODELING VAD SPACE	133
10.5. EMOTION PREDICTION IN THE PROPOSED FRAMEWORK	136
10.5.1. Multimodal Data Processing	137
10.5.1.1. Prediction of Emotion from a Visual Cue	138
10.5.1.2. Prediction of Emotion from Physiological Cue	139
10.5.2. Ground Truth Data	139
10.5.3. Emotion Prediction	140
10.6. FEATURE SELECTION AND CLASSIFICATION	140
10.7. RESULTS AND DISCUSSIONS	141
CONCLUSION	145
REFERENCES	146
SUBJECT INDEX	36:

FOREWORD

Affective Computing is a new area aiming to create intelligent computers that recognize, understand, and process human emotions. Affective Computing is an interdisciplinary area that encompasses a variety of disciplines, such as computer science, psychology, and cognitive science, among others. Emotion may be communicated in various ways, including gestures, postures, and facial expressions, as well as physiological signs, including brain activity, heart rate, muscle activity, blood pressure, and skin temperature. Humans can perceive emotion through facial expressions in general. However, not all emotions, particularly complex ones such as pride, love, mellowness, and sorrow, can be identified only through facial expressions. Physiological signals can therefore be utilized to represent complex emotions effectively.

This book aims to provide the audience with a basic understanding of Affective Computing and its application in many research fields. This state-of-the-art review of existing emotion theory and modeling approaches will help the readers explore various aspects of Affective Computing. By the end of the book, I hope that the readers will be able to understand emotion recognition methods based on audio, video, and physiological signals. Moreover, they will learn the fusion framework and familiarity to implement for emotion recognition.

Shitala Prasad
Institute for Infocomm Research
A*Star
Singapore

PREFACE

Affective Computing is an emerging field with the prime focus on developing intelligent systems that can perceive, interpret, process human emotions and act accordingly. Affective Computing incorporates interdisciplinary research areas like Computer Science, Psychology, Cognitive Science, Machine Learning, *etc.* Machines must perceive and interpret emotions in real-time and act accordingly for intelligent communication with human beings. Emotion plays a significant role in communication and can be expressed through many ways, like facial or auditory expression, gesture or sign language, *etc.* Brain activity, heart rate, muscular activity, blood pressure, and skin temperature are a few examples of physiological signals. It plays a crucial role in affect recognition compared to other emotion modalities. Humans perceive emotion primarily through facial expressions; yet, complex emotions such as pride, love, mellowness, and sorrow cannot be identified just by facial expressions. Physiological signals can thus be employed to recognize complex emotions.

The objective of this book is mainly three-fold: (1) Provide in-depth knowledge about affective Computing, affect information representation, models, and theories of emotions. (2) Emotion recognition from different affective modalities, such as audio, facial expression, and physiological signals, and (2) Multimodal fusion framework for emotion recognition in three-dimensional Valence, Arousal, and Dominance space.

Human emotions can be captured from various modalities, such as speech, facial expressions, physiological signals, *etc.* These modalities provide critical information that may be utilised to infer a user's emotional state. The primary emotions can be captured easily by facial and vocal expressions. However, facial expressions or audio information cannot detect complex emotions. Therefore, an efficient emotion model is required to predict complex emotions. The dimensional model of emotion can effectively model and recognize complex emotions.

Most emotion recognition work is based on facial and vocal expressions. However, the existing literature completely lacks emotion modeling in a continuous space. This book contributes in this direction by proposing an emotion model to predict a large number (more than fifteen) of complex emotions in a three-dimensional continuous space. We have implemented various systems to recognize emotion from speech, facial expression, physiological signals, and multimodal fusion of the above modalities. Our emphasis is on emotion modeling in a continuous space. Emotion prediction from physiological signals as complex emotions is better captured by physiological signals rather than facial or vocal expressions. The main contributions of this book can be summarized as follows:

1. This book presents a state-of-the-art review of Affective Computing and its application in various areas like gaming, medicine, virtual reality, *etc.*
2. A detailed review of multimodal fusion techniques is presented to assimilate multiple modalities to accomplish multimodal fusion tasks. The fusion methods are provided from the perspective of the requirement of multimodal fusions, the level of information fusion, and their applications in various domains, as reported in the literature. Moreover, significant challenges in multimodal fusions are also highlighted. Further, we present the evaluation measures for evaluating multimodal fusion techniques.
3. The significant contribution of this book is the three-dimensional emotion model based on valence, arousal, and dominance. The emotion prediction in three-dimensional space based on valence, arousal, and dominance is also presented.

CONSENT FOR PUBLICATION

Not applicable.

CONFLICT OF INTEREST

The author declares no conflict of interest, financial or otherwise.

Gyanendra K. Verma
Department of Information Technology
National Institute of Technology Raipur
Chhattisgarh
India

ACKNOWLEDGEMENTS

I acknowledge the guidance of Prof. U. S. Tiwary, IIIT Allahabad India, who motivated me to choose Affective Computing as a research topic many years ago. Some of his insights are still present in this book. I want to thank Dr. Shitala Prasad, Scientist, Institute for Infocomm Research, A*Star, Singapore, for his valuable suggestions. And last but not least, Ms. Humaira Hashmi, Editorial Manager Publications, Bentham Books, for extending her kind cooperation to complete this book project.

Affective Computing

Abstract: With the invention of high-power computing systems, machines are expected to show intelligence at par with human beings. A machine must be able to analyze and interpret emotions to demonstrate intelligent behavior. Affective computing not only helps computers to improve performance intelligently but also helps in decision-making. This chapter introduces affective computing and related issues that influence emotions. This study also provides an overview of human-computer interaction (HCI) and the possible use of different modalities for HCI. Further, challenges in affective computing are also discussed, along with the application of affective computing in various areas.

Keywords: Arousal, DEAP database, Dominance, EEG, Multiresolution analysis, Support vector machine, Valence.

1.1. INTRODUCTION

The cognitive, affective, and emotional information is crucial in HCI to improve user-computer connection [1]. It significantly enhances the learning environment. Emotion recognition is crucial since it has several applications in HCI and Human-Robot Interaction (HRI) [2] and many other new fields. Affective computing is a hot topic in the field of human-computer interaction. “Affective Computing is the research and development of systems and technologies that can identify, understand, process, and imitate human emotions,” according to the definition.

Affective computing is an interdisciplinary area that encompasses a variety of disciplines, such as computer science, psychology, and cognitive science, among others. Emotions can be exhibited in various ways, such as gestures, postures, facial expressions, and physiological signs, including brain activity, heart rate, muscular activity, blood pressure, and skin temperature [1].

People generally perceive emotion through facial expressions; nevertheless, complex emotions such as pride, gorgeousness, mellowness, and sadness cannot be identified through facial expressions [3]. Physiological signals can therefore be utilized to represent complicated effects.

1.2. WHAT IS EMOTION?

“Everyone knows what an emotion is until asked to give a definition.” [4].

Although emotion is prevalent in human communication, the term has no universally agreed meaning. Kleinginna and Kleinginna [5], on the other hand, gave the following definition of emotion:

“Emotion is a complex set of interactions between subjective and objective factors mediated by neural/hormonal systems that can:

1. Generate compelling experiences such as feelings of arousal, pleasure/displeasure;
2. Generate cognitive processes such as emotionally relevant perceptual effects, appraisals, and labeling processes;
3. Activate widespread physiological adjustments to arousing conditions; and
4. Lead to behavior that is often, but not always, expressive.”

1.2.1. Affective Human-Computer Interaction

The researchers described two ways to analyze emotion. The first method divides emotions into joy, fun, love, surprise, grief, *etc.* Another option is to display emotion on a multidimensional or continuous scale. Valence, arousal, and dominance are the three most prevalent aspects. How does a valence scale determine how happy or sad a person is? The arousal scale assesses how relaxed, bored, aroused, or thrilled [6]. The dominance scale depicts submissive (in control) or dominant (empowered) behavior. Emotion identification from facial expressions and voice signals is part of affective HCI. As a result, we will concentrate on the first two modalities, particularly concerning emotion perception. One of the essential needs of MMHCI is that multisensory data be processed individually before being merged.

A multi-modal system may be used in case of insufficient or noisy data. The system may use complementary information from other modalities if one modality's information is absent. If one modality fails to make a decision, the other must do so. Multi-modal HCI (MMHCI) incorporates several domains, such as Artificial Intelligence, Computer Vision, Psychology and others, according to Jaimes A. *et al.* [7]. People communicate frequently using facial expressions, bodily movement, sign language, and other non-verbal communication techniques [8].

Audio and video modalities are commonly employed in man-machine interaction; hence they are vital for HCI. At the feature or choice level, MMHCI focuses on merging several modalities of emotion. Probabilistic graphical models such as the Hidden Markov Model (HMM) and Bayesian Networks are beneficial, according to the study [9]. As a result of its ability to deal with missing values *via* probabilistic inference, Bayesian networks are widely used for data fusion. Vision methods are another option that may be employed for MMHCI [9]. The vision techniques categorize using a human-centered approach and decide how people may engage with the system.

1.3. BACKGROUND

Most emotion recognition research focuses on facial expression and voice emotion [10, 11, 12, 13]. Our book contributed to this approach by presenting an emotion model to predict many complicated emotions in a three-dimensional continuous space, lacking in the previous literature [14]. Even though we have created systems that identify emotion from speech, facial expression, physiological data, and multi-modal fusion of the modalities mentioned above, we focus on emotion modeling in a continuous space and emotion prediction using multi-modal cues.

People usually gather information from various sensory modalities, such as vision (sight), audition (hearing), tactile stimulation (touch), olfaction (smell), and gustation (taste). Then, this information is processed by integrating it into a single cohesive stream of information to communicate with others. In order to integrate numerous complementary and supplemental information, the human brain receives information from multiple communication modalities (such as reading text).

Multi-modal information fusion can be employed in effective systems to integrate related information from different modalities/cues to improve performance [15] and decrease ambiguity in decision-making by reducing data categorization uncertainty. Multi-modal information fusion is necessary for many applications where information from a single modality is inadequate and may contain noise or be insufficient to make conclusions. Consider a visual surveillance system where an object is monitored using visual information. If the object gets occluded, the surveillance system will have no way of tracking it.

Consider a surveillance system that takes information from two modalities (audio and visual information). The object can be tracked even if one of the modalities is unavailable; the system can process the information obtained from other modalities.

Affective Information Representation

Abstract: This chapter presents a brief overview of Affective computing and a formal definition of emotion given by various researchers. Human-computer interaction aims to enhance communication between man and machine so that machines can acquire, analyze, interpret and act on par with human beings. At the same time, Affective human-computer interaction focuses on enhancing communication between man and machines using affective information. Moreover, this chapter deals with Human emotional expression and perception through various modalities such as speech, facial expressions, physiological signals, *etc.* It also detailed the overview of Action Units and Techniques for classifying facial expressions as reported in the literature.

Keywords: Action units, Affective computing, Affective HCI, Emotion expression, Facial expression, HCI.

2.1. INTRODUCTION

A review of multimodal emotional information extraction and processing is presented in this chapter. Affective computing may be thought of as an issue of automatic emotion perception for improved human-machine connection. It entails the detection and interpretation of human emotion as well as the prediction of the user's mental state. It also entails the examination of a person's emotional data in order to determine his or her mental state. Facial emotion identification is explored in the chapter, along with the principles of facial expression and emotion modeling. With the system's limitations, facial expression representation is also available.

2.2. AFFECTIVE COMPUTING AND EMOTION

Affective computing is an interdisciplinary area that encompasses a variety of disciplines, such as computer science, psychology, and cognitive science, among others. Emotion may be exhibited in various ways, including gestures, postures, facial expressions, and physiological signs, heart rate, muscular activity, blood pressure, and skin temperature [1].

According to the definition, “Affective Computing is the research and development of systems and technologies that can identify, understand, process, and imitate human emotions”.

People generally perceive emotion through facial expressions; nevertheless, complex emotions such as pride, gorgeousness, mellowness, and sadness cannot be identified through facial expressions [2]. Physiological signals can therefore be utilized to represent complicated effects.

2.2.1. Affective Human-Computer Interaction

The researchers described two ways to analyze emotion. The first method divides emotions into categories, such as joy, fun, love, surprise, grief, *etc.* Another option is to display emotion on a multidimensional or continuous scale. Valence, arousal and dominance are the three most prevalent aspects. How does a valence scale determine how happy or sad a person is? The arousal scale assesses how relaxed, bored, aroused, or thrilled you are [3]. The dominance scale depicts submissive (in control) or dominant (empowered) behavior. Emotion identification from facial expressions and voice signals is part of affective HCI. As a result, we will concentrate on the first two modalities, particularly concerning emotion perception. One of the essential needs of MMHCI is that multisensory data be processed individually before being merged.

Multimodal HCI (MMHCI) incorporates several domains, such as Artificial Intelligence, Computer Vision, Psychology, and others, according to Jaimes A. *et al.* [4]. When people connect, they frequently employ both verbal and nonverbal communication. Non-verbal communication includes facial expressions, bodily movement, sign language, and other non-verbal techniques [5].

The most commonly used modalities in HCI are audio and video; hence, they are vital for HCI. Multimodal HCI focuses on merging several modalities of emotion at the feature or decision level. Probabilistic graphical models such as the Hidden Markov Model (HMM) and Bayesian Networks are beneficial, according to [6]. As a result of its ability to deal with missing values, Bayesian networks are commonly utilized for data fusion through probabilistic inference. Vision methods are another option that may be employed for MMHCI [6]. The vision techniques categorize using a human-centered approach and decide how people may engage with the system.

2.2.2. Human Emotion Expression and Perception

“Everyone knows what an emotion is, until asked to give a definition.” [7].

Although emotion is prevalent in human communication, the term has no universally agreed meaning. (Kleinginna and Kleinginna) [8].

On the other hand, it gave the following definition of emotion:

1. “Emotion is a complex set of interactions between subjective and objective factors mediated by neural/hormonal systems that can:
2. Generate compelling experiences, such as feelings of arousal, pleasure/displeasure; b) Generate cognitive processes, such as emotionally relevant perceptual effects, appraisals, and labeling processes;
3. Activate widespread physiological adjustments to arousing conditions; and
4. Lead to behavior that is often, but not always, expressive.”

Automatic Human Emotion Recognition captures and extracts information from numerous emotional modalities. We have a wide range of sensors and devices to gather voice data, visual signals, linguistic contents, and physiological signals, among other things, in order to capture emotional information. Although emotional information may be expressed in various ways, Fig. (2.1) depicts the most commonly utilized signals reported in the literature.

2.2.2.1. Facial Expressions

The majority of the study has focused on detecting emotion through facial expressions. One of the most dependable and natural ways to communicate emotion is through facial expression. We may quickly detect another person's enjoyment, grief, disagreement, and intentions during the conversation by facial expression. Another advantage of facial expression is that anybody may show it, regardless of age or gender. As a result, affective computing's primary source/channel is facial expression.

2.2.2.2. Audio and Speech

For spoken communication, audio is the most common channel. A speech reflects the style of communication by conveying emotion through linguistic and paralinguistic signals. In any event, a person's affective state may be deduced directly from speech. Emotion theory shows that a physiological response accompanies the emotional state. When we are walking through a jungle, we are

Models and Theory of Emotion

Abstract: This chapter presents a state-of-the-art review of existing emotion theory, modeling approaches, and affective information extraction and processing methods. The basic theory of emotions deals with Darwin's evolutionary theory, Schechter's theory of emotion, and James–Lange's theory. These theories are fundamental building blocks of Affective Computing research. Emotion modeling approaches can be categorized into categorical, appraisal, and dimensional models. Noticeable contributions to Affect recognition systems in terms of modality, database, and dimensionality are also discussed in this chapter.

Keywords: Appraisal model, Categorical model, Dimensional model, Emotion modeling, Emotion theory.

3.1. INTRODUCTION

Modeling emotion is essential for a better understanding of emotions. Efficient modeling of emotions is still challenging due to the involvement of various emotional modalities. As each modality has a different pattern, emotion modeling depends upon the type of input signals. Emotion theories and models are the basis for gaining in-depth knowledge about the induction of emotions. Researchers have proposed many emotion theories; among them, the James-Lange theory, Canon-Bard theory, and Schachter-Singer theory of emotion contributed significantly.

The emotional state of a human being at a particular moment is the combination of the user's physiological, psychological, and subjective experience [1]. The appraisal experience depends on various parameters, such as growing environment, background, and culture. Thus, people feel different experiences of similar phenomena or events.

3.2. EMOTION THEORY

The three primary emotion models are the category, appraisal-based, and dimensional models [2].

The category model is concerned with universally recognized basic emotions. The number of these basic emotions is modest, yet they are all linked to our brains [3]. The appraisal-based approach deals with modelling emotions as a physiological response to stimuli or events that leads to the emergence of emotion and associated action. On a continuous or discrete scale, dimensional methods describe emotion through some independent dimensions. Fig. (3.1) depicts various emotion theories.

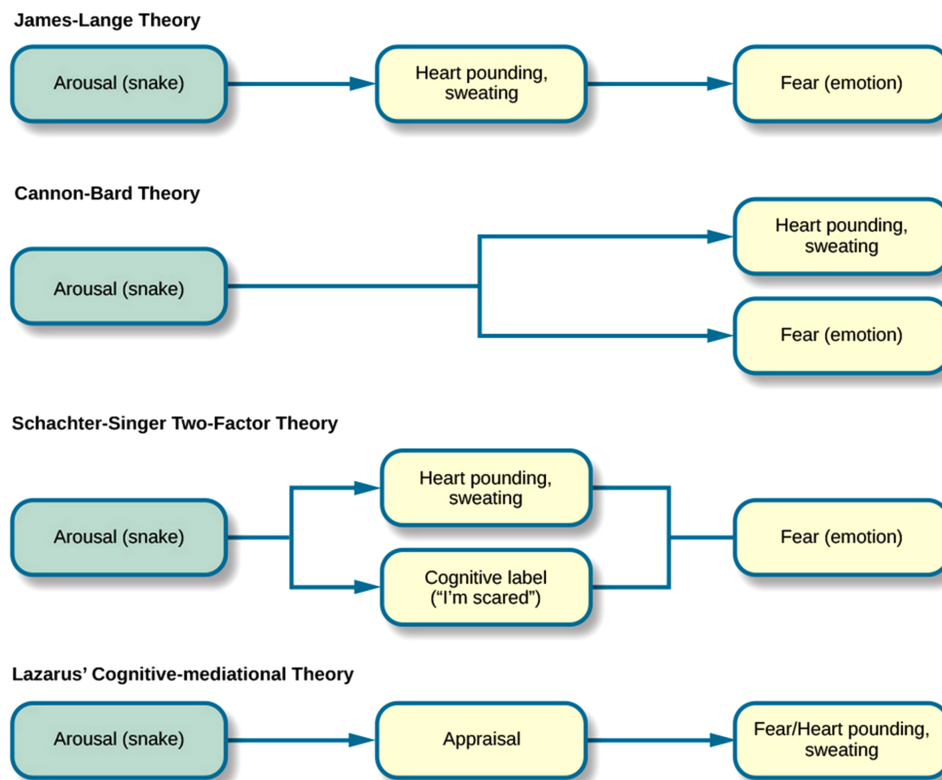


Fig.(3.1). Major emotion theories.

3.2.1. Categorical Approach

Prof. P. Ekman made significant contributions in the area of Emotion Recognition. He was the first to introduce six basic emotions visible in facial expressions [4]. Fig. (3.2) illustrates the six primary emotions proposed by Ekman. In 1980, Robert Plutchik introduced the “Wheel of emotion,” a new notion of emotion [5]. Fig. (3.2) classified a few feelings as core emotions: joy, sorrow, anger, fear, trust, contempt, surprise, and anticipation. In 2001, Parrot offered a classification of emotion, which Robert Plutchik followed.



Fig. (3.2). Ekman's six universal emotions.

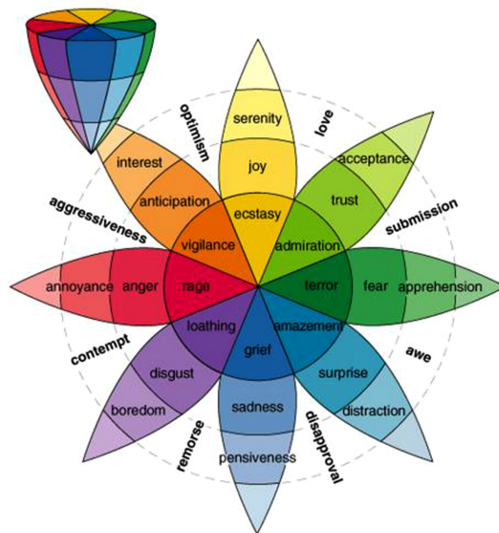


Fig. (3.3). Robert Plutchik's Wheel of emotion [8].

The idea of Parrot was to divide emotions into three categories: primary, secondary, and tertiary [6]. There are six significant emotions, twenty-five secondary emotions, and more than tertiary ones in Parrots' hypothesis. Among the different theories of emotions published in the literature, Plutchik and Conte's [7] study is the most noteworthy. The "Psycho-evolutionary Theory of Emotions" was created by them. Robert Plutchik's wheel of emotion [8] is shown in Figure (3.3)

3.2.2. Evolutionary Theory of Emotion by Darwin

Ekman [9] defined the six primary emotions as "happy," "sad," "angry," "fear," "surprise," and "disgust." These six fundamental emotions are based on Darwin's

Affective Information Extraction, Processing and Evaluation

Abstract: This chapter presents a state-of-the-art review of existing affective information extraction and processing approaches. Various evaluation criteria, such as Evaluation matrices like ROC, F1 measure, Mean Square Error, Mean Average Error, Threshold criteria, and Performance criteria are also reported in this chapter.

Keywords: Evaluation measures, F1 measures, Information extraction, ROC curve.

4.1. INTRODUCTION

Information extraction from multimodal cues is usually used as complementary. When getting the information from one cue is terminated, the system may shift to other traits and continue its functioning. There are three primary sources of affective information: speech, images or videos, and text. The first two cues, *i.e.*, speech and images, are widely used for affective information processing. However, text-based affective information analysis targets specific applications and comes under Natural Language Processing (NLP) domain [1]. Where information processing is concerned, the time-series data is easy to process and extract. Usually, the audio and image information is synchronized in a time-dependent manner.

4.2. AFFECTIVE INFORMATION EXTRACTION AND PROCESSING

4.2.1. Information Extraction from Audio

Diverse speech aspects, such as mood, speaker audio, and both, are used to represent different speech information. As a result, experts are interested in learning more about the speech characteristics of various emotions. Vocal tract, prosodic, and excitation source characteristics are the three types of speech features. Short segments of voice signals are used to extract Vocal Tract characteristics. The energy distribution for a variety of speech frequencies is

represented by these properties. Various vocal tract elements and their combinations are employed by various studies for emotion identification. T. Long [2] employed a mixture of Perceptual Linear Prediction (PLP), MFCC, and LPCC to distinguish emotions, such as angry, happy, sad, bored, and neutral, using the log frequency power coefficient (LFPC) vocal tract feature. As pitch has a higher discriminating power than other prosodic variables, it is the most extensively employed prosodic feature for emotion identification. Aside from pitch, log energy is a popular metric for analysing speaking styles and emotions.

4.2.2. Information Extraction from Video

The appearance of changes owing to lighting and position fluctuations complicates geometric feature-based face analysis. Hence, spatio-temporal characteristics can be utilised to identify minor changes in a face. The interest spots in picture sequences are detected using Dollar's approach [3]. This approach was created to identify minor changes in the spatial and temporal domains in which the human action recognition community is most interested. Multiresolution techniques may also be used to extract visual information from photographs.

4.2.3. Information Extraction from Physiological Signals

Physiological signals are vital for emotion, according to several emotion theories. Ekman established that a certain emotion may be linked to a specific physiological pattern. The frequency and amplitude of an EEG signal can be used to characterize it.

To split the physiological signal into distinct frequency bands, a bandpass filter can be utilised. Discrete Wavelet Transform is another approach for decomposing EEG signals into distinct frequency bands (DWT). Galvanic Skin Response (GSR), respiration amplitude, electrocardiogram (ECG), electromyograms (EMG), electrooculogram (EOG), Electro dermal Activity (EDA), Galvanic Skin Response (GSR), Skin Conductance Response (SCR), and skin temperature are examples of peripheral biosignals. EDA and GSR are skin conductance measurements that are extensively utilised for automated emotion identification. GSR is a fairly reliable physiological marker of human arousal, according to J. Kim [4].

4.3. STUDIES ON AFFECT INFORMATION PROCESSING

People use various modes of communication in day-to-day interaction with others. We can broadly categorize these communication modes into verbal and non-verbal communication. Verbal communication involves speech and audio; however, non-verbal communication uses facial expressions, body gestures, and sign language. Both communication modes are vital and play a significant role in communication among human beings. Unfortunately, we lack a better human-computer interface to use these vital communication channels. Facial expressions are essential to recognize sentiments in communications. Koelstra *et al.* [5] discovered substantial changes in N400 ERP responses when placing relevant and irrelevant tags on short videos. They started by extracting features from audio and video channels. Then, the correlation between audio and visual properties was investigated using the method described above. After that, a hidden Markov model is utilized to characterize statistical dependency across time segments and discover the features in the altered domain's fundamental temporal structure. Extensive system testing is used to assess the resilience of our suggested solution.

Emotion identification can be unimodal or multimodal. The unimodal-based method uses a single modality to recognize the emotion. In contrast, the multimodal technique collects emotional information from several inputs, such as audio, video, image, physiological signals, *etc.* We have dealt with multimodal emotion recognition in this study. Several researchers also used this multimodal strategy. S. Koelstra *et al.* [5] suggested a facial expression and EEG signal fusion system for emotion identification using multimodal fusion.

Mamalis A. Nikolaou [6] demonstrated a multimodal emotion identification system that combines facial expressions, shoulder gestures, and aural signals. They used two-dimensional spaces of valence and arousal to map multimodal emotions on a continuous scale. Further, an associative fusion framework was presented using Support Vector Regression (SVR) and Long Short Term Memory neural networks (BLSTM-NNs).

Y. Wang *et al.* [7] used the kernel approach to examine multimodal information extraction and analysis. For modeling the nonlinear interaction between two multidimensional variables, they used Kernel cross-modal factor analysis. They have also developed a method for determining the best transformations to describe patterns.

M. Paleri *et al.* [8] published a paper on feature selection for automated audio-visual person-independent emotion identification. They employed a neural network to compare the performance of different characteristics in an emotion identification system.

Multimodal Affective Information Fusion

Abstract: The multimodal information can be assimilated at three levels 1) early fusion, 2) intermediate fusion, and 3) late fusion. Early fusion can be performed at the sensor or signal level. Intermediate fusion can be at the feature level, and late fusion may be done at the decision level. Apart from that, some more fusion techniques are rank-based, adaptive, *etc.* This chapter provides an extensive review of studies based on fusion and reported noticeable work herewith. Eventually, we discussed the challenges associated with multimodal fusion.

Keywords: Decision fusion, Feature fusion, Multi-modal fusion, Sensor fusion.

5.1. INTRODUCTION

Affective information plays a significant role in emotion recognition. The information acquired from multiple sources or modalities, such as audio, video, and text, is known as Multimodal information. Multimodal information fusion is described as merging information from many sources/modalities to produce superior performance than a single source/modality [1]. There have been several fusion categories mentioned in the literature. The fusion might happen at the signal, feature, or decision level. Early fusion occurs when the fusion is accomplished at the signal or feature level. Late fusion refers to fusion that occurs after a choice has been made. It is not required for distinct modalities to give complementary information in the fusion process, according to P. K. Atrey *et al.* [2], hence it is essential to know which modalities contribute the most. The appropriate number of modalities in the fusion process is a critical component.

5.2. MULTIMODAL INFORMATION FUSION

The early fusion integrates the information acquired from different modalities before applying learning models. Late fusion (also known as decision-level fusion) integrates the information acquired from the output of different algorithms or models. According to L. Hoste [3], Late fusion is based on the semantic information fusion obtained from different modalities. Each modality should be processed first, then integrated at the end to handle multimodal data. Several st-

udies [4 - 6] discussed fusion architectures and multimodal data processing. The data in multimodal processing is not necessarily mutually independent and, therefore, cannot be merged in a context-free fashion. However, the information must be processed in a combined space using a context-dependent model. The critical challenges in multimodal fusion are different feature formats. The dimensionality of joint feature space with temporal synchronization is another issue [4].

The two most significant aspects of multimodal information fusion are i) level of fusion and ii) the type of fusion. The fusion process must be synchronized. Time complexity and cost-effectiveness are the other performance factors. The fusion of two or more modalities must be done methodically. The primary difficulties in multimodal fusion are the number of modalities and information synchronization. The primary difficulties are a correlation between information collected from multiple modalities and the level at which the information is fused [7, 8]. During fusion, many modalities may not necessarily give supplementary information. As a result, it is critical to comprehend the contributions of each modality. The prime role of Multimodal information fusion is to combine data from different modalities/cues to eliminate ambiguity and uncertainty.

Information can be derived from various sources/modalities, such as text, images, and speech. The information can have low and high-level characteristics depending on the input source. The features are fused at either a low (feature fusion) or a high (decision level) level. The fusion techniques and associated works documented in the literature are also presented in this chapter.

5.2.1. Early Fusion

Information can be fused early, *i.e.*, at the sensor or signal level. A three-dimensional image, for example, can be created by fusing two or more two-dimensional images [9]. An example of early fusion is audio-visual information fusion, which integrates the audio and video information into a single feature vector. Dimensional reduction methods like principal component analysis or linear discriminate analysis can be used for dimensionality reduction. The two modalities are combined even before classification using Adaptive fusion. In Adaptive fusion, the weights are assigned to each modality, which is one of the essential advantages of early fusion. When we wish to give a modality greater weight, adaptive fusion comes in handy.

Numerous approaches have been reported in the literature to execute fusion, including Bayesian inference, Dempster–Shafer fusion, Maximum Entropy model, and Nave Bay's algorithms. Pitsikalis *et al.* [10] suggested a Bayesian

inference-based technique integrating audio-visual information. They calculated the joint probability of integrated MFCC and texture analysis characteristics.

Mena and Malpica employed a Dempster–Shafer fusion technique for color picture segmentation [11]. For semantic multimedia indexing, Magalhaes and Ruger [12] employed the maximum entropy approach. For image retrieval, they integrated text and picture characteristics.

5.2.2. Intermediate Fusion

The drawback of the early fusion technique is its inability to deal with erroneous data that avoids explicit modeling of the various modalities. Early fusion methods suffer from relative dependability, fluctuations, and asynchrony. Such issue may be resolved by comparing the time instance feature to the time scale dimension of the relevant stream. As a result, by comparing previously seen data instances with current data transmitted *via* a working observation channel, one may create a statistical forecast with a derivable probability value for erroneous instances owing to sensor failures, *etc.*

Probabilistic graphical models are best suited for fusing numerous sources of data in this context [13]. Probabilistic inference also handles noisy features and missing feature values. A hierarchical HMM to identify facial expressions was presented by Cohen *et al.* [14]. The capability to fuse multiple sources of information of dynamic Bayesian networks and HMM variations was demonstrated by Minsky [15]. Carpenter R [16]. proposed a fusion approach to detect office activity and events in video utilizing audio and video signals.

5.2.3. Late Fusion

A multimodal system combines many modalities in order to conclude. It demands a standard framework for representing shared meaning across all modalities and a well-defined method for information assimilation [4]. Late fusion models typically employ distinct classifiers for each stream that are trained separately. The outcome is derived by fusing the outcome of individual classifiers. Only at the integration stage is the correspondence between the channels recognized. Late fusion has several clear advantages.

The inputs may be recognized independently, they do not have to coincide. The late fusion system employs classifiers that can be trained on a single data set but are scalable in terms of the number of modes and vocabulary. As with audio-visual recognition, we need to find a decent heuristic for extracting features from

Multimodal Fusion Framework and Multiresolution Analysis

Abstract: This chapter presents a multi-modal fusion framework for emotion recognition using multiresolution analysis. The proposed framework consists of three significant steps: (1) feature extraction and selection, (2) feature level fusion, and (3) mapping of emotions in three-dimensional VAD space. The proposed framework considers subject-independent features that can incorporate many more emotions. It is possible to handle many channel features, especially synchronous EEG channels and feature-level fusion works. This framework of representing emotions in 3D space can be extended for mapping emotion in three-dimensional spaces with three specific coordinates for a particular emotion. In addition to the fusion framework, we have explained multiresolution approaches, such as wavelet and curvelet transform, to classify and predict emotions.

Keywords: Curvelet transform, Emotion recognition, Wavelet transform.

6.1. INTRODUCTION

Multimodal fusion combines multiple cues that may act as complementary information to improve the system's performance. Several fusion approaches are reported in the literature; nonetheless, early, middle and late fusion are three primary categories. Before feeding into the learning phase, the features gathered from diverse modalities must be integrated into a single representation in an early fusion. Intermediate fusion can deal with insufficient data and asynchrony across distinct modalities. Decision-level fusion deals with semantic information as the decision is made by considering the outcome of different modalities after feature extraction. One of the essential criteria for multimodal data processing is that the data be processed individually before joining.

6.2. THE BENEFITS OF MULTIMODAL FEATURES

The challenging issues in multimodal fusion are 1) data types of modalities, 2) the synchronization of different types of modalities, and 3) the level of the fusion [1, 2]. The choice of fusion level may be easy for a similar type of data. For example,

if the two modalities are temporal, the fusion may be easy; however, if one is temporal and the other is spatial, fusion becomes challenging. The contribution of each modality is unique, and each fusion process does not need to enhance the system's performance.

Multiple modalities, such as text, picture, audio, and physiological inputs, can all have features extracted individually. The multimodal element not only adds to the information available but also has the potential to improve the system's performance. The following are some of the reasons for not relying on a single method of features:

6.2.1. Noise In Sensed Data

Sensor, channel, and modality-specific noise are the three forms of noise in sensed data. Sensor noise is the noise produced by the sensor. Each pixel of a camera sensor, for example, is made up of one or more light-sensitive photodiodes that convert incoming light into an electrical signal. The color value of the final image pixel represents the signal. Even if the same pixel were exposed to the same quantity of light numerous times, the resultant color value would not be equal. However, it would have a slight variance known as “noise.”

Channel noise, on the other hand, is the result of the data transmission or medium deteriorating. Under slightly varied lighting conditions, for example, the same HCI modality may change. Person identification is perhaps the most well-known example. Under varied lighting circumstances, the same face looks different from two separate faces recorded under the same illumination conditions. Finally, modality-specific noise is noise induced by a disagreement between the acquired data and the standard interpretation of the modality.

6.2.2. Non-Universality

A system may not be able to get valuable data from only one modality. For example, complex emotions, such as pride, joy, excitement, sorrow, *etc.*, cannot be identified by facial expressions only. Thus we must rely on other methods, such as physiological signals, to recognize complex emotions. Similarly, iris recognition biometrics may fail due to different eye conditions like long eyelashes, sloping eyelids, or certain eye pathologies issues. On the other hand, a face recognition system may nevertheless be a valuable biometric modality. While no one modality is ideal, combining them should provide more excellent user coverage, enhancing accessibility, particularly for the impaired.

6.2.3. Complementary Information

The information gained from the other modality can be utilized as a supplement. A single modality-based algorithm may fail if the input signal is lost or corrupted. A unimodal system may stop functioning in case of inputs from one modality interrupt. However, a multimodal system can continuously perform by taking inputs from other modalities. An object-tracking method based on visual modality works perfectly in the usual scenario. However, if the object occluded, it will stop tracking; in this instance, the voice modality may provide complementary information.

6.3. FEATURE LEVEL FUSION

Multimodal information fusion is the job of combining corresponding data from different modalities/cues in order to eliminate ambiguity and uncertainty. Depending on the applications, information can be acquired from various sources/modalities, such as text, pictures, and speech. The low and high-level features extracted from various modalities can be fused at either the feature level or a decision level.

Sensor level fusion combines raw data or data generated from several sources from sensory (raw) data. The finest example of sensor-level fusion is creating a 3D picture from two 2D images. The type of sensors and information sources are two essential concerns in sensor-level fusion. Other important characteristics for sensor level fusion are as follows [3]:

- Sensors' computational capability
- Topology, communication structure, computing resources, and
- System goals and optimization
- Improved detection, tracking, and identification are some of the benefits of sensor-level data fusion.
- Improved situational awareness and assessment
- Increased sturdiness
- Coverage that is both spatially and temporally extensive
- Reduced communication and computing costs, as well as a faster reaction time.

Feature level fusion is achieved by extracting features from several modalities/sources individually and then combining them after normalization. The benefit of feature-level fusion is that it can find the correlation between feature vectors, which can help the system perform better. The following parameters affect feature-level fusion.

Emotion Recognition From Facial Expression In A Noisy Environment

Abstract: This study presents emotion recognition from facial expressions in a noisy environment. The challenges addressed in this study are noise in the images and illumination changes. Wavelets have been extensively used for noise reduction; therefore, we have applied wavelet and curvelet analysis from noisy images. The experiments are performed with different values of Gaussian noise (mean: 0.01, 0.03) and (variance: 0.01, 0.03). Similarly, for experimentation with illumination changes, we have considered different dynamic ranges (0.1, 0.9). Three benchmark databases, Cohn-Kanade, JAFFE, and In-house, are used for all experimentation. The five best machine learning algorithms are used for classification purposes. Experimental results show that SVM and MLP classifiers with wavelet and curvelet-based coefficients yield better results for emotion recognition. We can conclude that Wavelet coefficients-based features perform well, especially in the presence of Gaussian noise and illumination changes for facial expression recognition.

Keywords: Cohn-Kanade, Curvelet transform, Facial expression recognition, JAFFE, MLP, SVM, Wavelet transform.

7.1. INTRODUCTION

Emotions are essential for machines to make an intelligent decision. We are witnessing exponential growth in computing power, however, lacking robust algorithms that can enhance the intelligence of the machines. Emotions play a significant role in an intelligent behavior by machines at par with human beings [1, 2, 3, 4]. Emotions can be exhibited through various modes, such as facial expression, auditory expression, physiological expression, gesture, body language, sign language, *etc.* Facial expression is the most widely used modality among the above modalities due to quick presentation and recognition. Moreover, facial expressions are more accessible to acquire, process, and analyze than other modalities.

Multiresolution Analysis (MRA) proved useful in various applications, including medical imaging, satellite imaging, biometrics, *etc.* Wavelet and Curvelet transforms are two classical algorithms used for MRA. In [5, 6], wavelet

transform-based multiresolution analysis was performed. However, the application of MRA in a noisy environment is relatively new for emotion recognition.

Some of the research [7] works are based on applying multiple modalities rather than a single modality to analyze emotions. Mansoorizadeh *et al.* [8] proposed a multimodal fusion framework for human emotion recognition. Some of the work is based on curvelet analysis [9, 10]. Lee and Shih [11] presented contourlet analysis with regularized discriminate analysis for facial emotion recognition. Shan *et al.* [12] presented a facial expression recognition system using local binary patterns. M. Yeasin *et al.* [13] proposed an approach for the measurement of levels of interest from video for facial expression recognition.

Generally, raw images are degraded due to various phenomena such as varying lighting conditions, environmental effects, high/ low brightness, contrast, *etc.* Thus, facial expression recognition became more challenging in a noisy environment. This chapter deals with the recognition of facial expressions in a noisy environment. We have experimented with three benchmark databases to prove the proposed algorithm's usefulness and robustness. Various types of noise have been added with different mean and variance values. Then, the multiresolution approaches were applied to extract the most prominent features of noisy facial expression images. Getting the desired accuracy of the system in a noisy environment is still a challenge [14]. This chapter proposes multiresolution approaches based on wavelet and curvelet analysis to improve emotion recognition performance.

7.2. THE CHALLENGES IN FACIAL EMOTION RECOGNITION

The major challenge in facial expression recognition is noise and illumination change. The performance of most of the algorithms degraded due to the presence of these two parameters: 1) the presence of noise in facial expression and 2) varying illumination. We have recreated a noisy database by adding Gaussian white noise with different mean and variance values. At the same time, we have also modified the dynamic range of the test images by 0.1 to 0.9. The sample database images having different noise and illumination changes are shown in Fig. (7.1).

The images may be degraded during acquisition due to environmental conditions such as lighting, illumination, handshake, *etc.* The noise may also be added due to the sensor error during raw data acquisition. Many systems cannot preprocess or remove such added noise in the images. Therefore, a robust system is required to handle and analyze such noisy data efficiently. In an image, the edge information

is degraded due to noise. It may happen even due to contrast reversals in some cases.

We have prepared the database for the experiments by adding Gaussian white noise in the images at different degrees. The values of mean and variance were kept in the range of 0.01- 0.03. For the first iteration of the experiment, the mean and variance values are set as 0.1 and 0.1. For the second test, the values are 0.03 and 0.03. Fig. (7.2) illustrates some sample images having additive noise. The research revealed that emotional signals could withstand noise-induced distortion. As all multiresolution approaches are inherently scale-invariant, we did not test for scale variation.

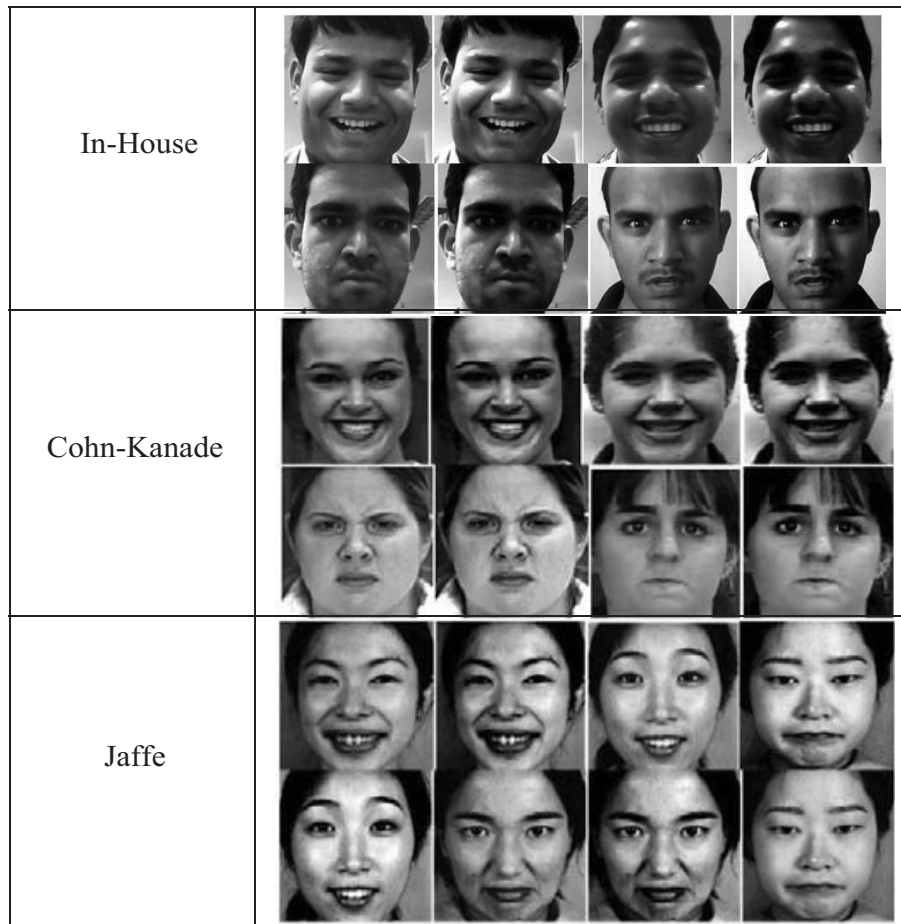


Fig. (7.1). Sample images with varying illumination conditions.

CHAPTER 8**Spontaneous Emotion Recognition From Audio-Visual Signals**

Abstract: This chapter introduces an emotion recognition system based on audio and video cues. For audio-based emotion recognition, we have explored various aspects of feature extraction and classification strategy and found that wavelet analysis is sound. We have shown comparative results for discriminating capabilities of various combinations of features using the Fisher Discriminant Analysis (FDA). Finally, we have combined the audio and video features using a feature-level fusion approach. All the experiments are performed with eINTERFACE and RML databases. Though we have applied multiple classifiers, SVM shows significantly improved performance with a single modality and fusion. The results obtained using fusion outperformed in contrast results based on a single modality of audio or video. We can conclude that fusion approaches are best as it is using complementary information from multiple modalities.

Keywords: Emotion recognition system, eINTERFACE, Feature fusion, Machine learning, RML database, Support vector machine.

8.1. INTRODUCTION

Most emotion recognition work is based on facial expression and vocal emotion [1, 2, 3, 4]. However, emotion modeling in a continuous space is completely lacking in the existing literature [5]; our book contributed in this direction by proposing an emotion model to predict a large number (more than fifteen) of complex emotions in a three-dimensional continuous space. Although we have implemented various systems to recognize emotion from speech, facial expression, physiological signals, and multimodal fusion of the above modalities, our emphasis is on emotion modeling in a continuous space.

This chapter investigates the use of audio-visual signals for spontaneous affect recognition. The proposed approach is based on Multiresolution analysis (MRA), which can analyze a signal at many resolutions. An MRA is the design process for the most practically important discrete wavelet transforms (DWT) and the reasoning for the Fast Wavelet Transform algorithm (FWT). MRA uses DWT to

extract the features. Different classifiers are used to classify the data, including the Support Vector Machine (SVM), Multilayer Perceptron (MLP), and K mean classifier. The MRA-based method in this book was created by building on prior work in emotion recognition using MRA [6]. Emotion recognition using audio and visual modalities is covered in this chapter.

The following is a breakdown of the chapter's structure. We review the many techniques for emotion recognition that have been reported in the literature. Then, using a development set of eINTERFACE [7] and RML [8] multimodal emotion databases, we present an emotion detection system based on auditory and visual cues and some of its main components. The use of multimodal fusion for emotion recognition is then investigated. Finally, experimental findings for audio and visual cues, as well as their combination utilizing feature level fusion, are shown on the eINTERFACE and RML databases.

8.2. RECOGNITION OF SPONTANEOUS EFFECTS

Affect recognition has been an emerging study subject within Human-Computer Interaction (HCI). Several affective states, including thinking, shame, and despair, are considered complex affective states and communicated *via* hundreds of distinct facial expressions. Researchers employed multimodal cues to identify emotion as a single cue may not correctly describe complex emotional states [9]. Several methods have been reported in the literature to extract characteristics from audio and visual cues.

8.3. THE DATABASE

Many academicians have been inspired to construct an emotion database due to recent breakthroughs in emotion identification.

- MIT [10],
- MMI [11],
- HUMAINE [12],
- VAM [13],
- SEMAINE [14],
- MAHNOB-HCI [15], and
- DEAP [16].

We used the eINTERFACE and RML audio-visual databases in this study. Fig. (8.1) depicts the facial expressions of various emotions for the eINTERFACE and RML databases. At the same time, Table 8.1 summarizes the database content for both databases.

8.3.1. eINTERFACE Database

The eINTERFACE05 is an audio and video emotion database that serves as a benchmark. Happy, Angry, Disgust, Sadness, Surprise, and Fear are the six primary emotions in the database. Forty-two people were involved in the collection of speech samples.



Fig. (8.1). Facial expressions for different emotions first row - eINTERFACE database, second row- RML database.

Table 8.1. Database content summary of eINTERFACE and RML database.

eINTERFACE	
Database type	Audiovisual
No. of subjects	44
Language	English
# Emotion	6 Universal emotions
#Video	1320
Image frame size	720*576
Frame rate	25
Audio sampling rate	48000
RML	
Database type	Audiovisual
No. of subjects	8
Language	6 (English, Mandarin, Urdu, Punjabi, Persian, Italian)

Multimodal Fusion Framework: Emotion Recognition From Physiological Signals

Abstract: This study presents a multimodal fusion framework for emotion recognition from physiological signals. In contrast to emotion recognition through facial expression, a large number of emotions can be recognized accurately through physiological signals. The DEAP database, a benchmark multimodal database with many collections of EEG and peripheral signals, is employed for experimentation. The proposed method takes into account those features that are subject-independent and can incorporate many more emotions. As it is possible to handle many channel features, especially synchronous EEG channels, feature-level fusion is applied in this study. The features extracted from EEG and peripheral signals include relative, logarithmic, and absolute power energy of Alpha, Beta, Gamma, Delta, and Theta. Experimental results demonstrate that physiological signals' Theta and Beta bands are the most significant contributor to the performance. On the other hand, SVM performs outstandingly.

Keywords: 3D emotion model, EEG, Facial expression, Fusion framework, Physiological signal.

9.1. INTRODUCTION

Physiological cues are essential in determining whether or not a person's behavior or emotional state has changed. Physiological signals are biosignals such as Galvanic Skin Response (GSR), Respiration amplitude, Electrocardiogram (ECG), Electro-myograms (EMG), Electrooculogram (EOG), Electro Dermal Activity (EDA), Galvanic Skin Response (GSR), Skin Conductance Response (SCR), and skin temperature [1]. EDA and GSR are skin conductance measurements extensively utilized for automated emotion identification. According to Kim J., GSR is a relatively reliable physiological marker of human arousal [2]. EEG, EMG, Skin conductance, BVP, and other physiological signals can be employed for emotion identification.

Electroencephalogram (EEG) is used in this study. EEG measures voltage changes in the brain's neurons induced by ionic current flows. For a brief length of time, 20–40 minutes, several electrodes put on the scalp capture the brain's spont-

aneous electrical activity. A brief introduction of various bio signals is given. A typical architecture of emotion recognition from physiological signals is illustrated in Fig. (9.1).

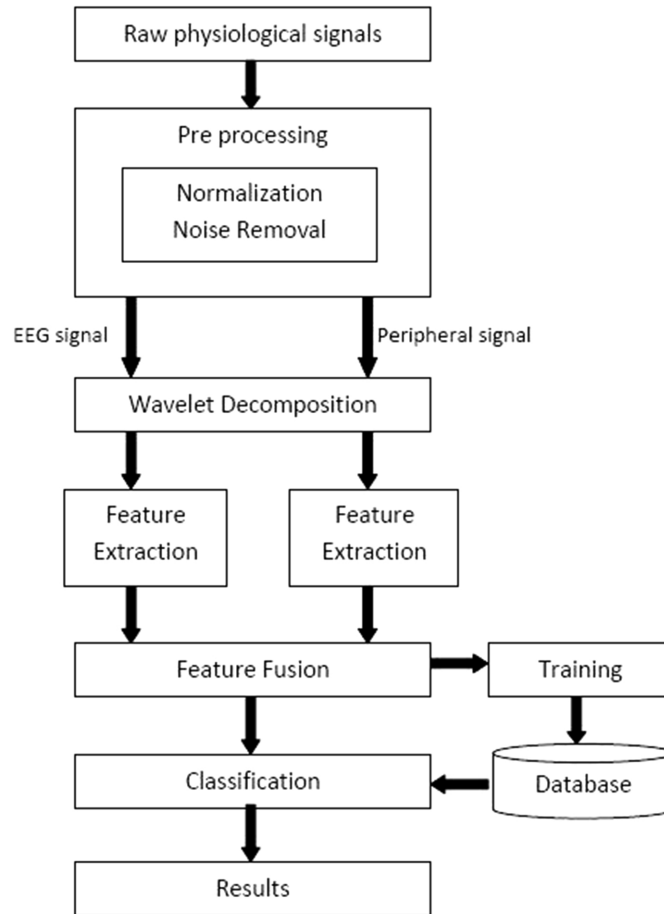


Fig. (9.1). Architecture of emotion recognition from physiological signal.

9.1.1. Electrical Brain Activity

EEG is a technique for detecting continuous electrical brain activity with amplitudes ranging from 1 to 200 microvolts. The 10-20 electrode placement approach is the most prevalent. Hundreds of electrodes (roughly) are implanted in the skull depending on its geometric dimensions.

9.1.2. Muscle Activity

Electromyography measures the electrical potentials that muscle fibers generate from muscular contraction (EMG). Ekman *et al.* [3] stated that facial EMG could give a sensitive and objective metric for emotion recognition as some facial muscles cannot be engaged voluntarily.

9.1.3. Skin Conductivity

The skin's electrical resistance is measured by skin conductance and regulated by the activity of the preparatory glands. The greater the skin conductance, the more active the preparation glands are.

9.1.4. Skin Temperature

According to McFarland [4], stimulating negative emotions causes skin temperature to drop, but calm, pleasant emotions cause skin temperature to rise. Emotion identification may also be made using heart rate and respiratory rate. Different physiological signals may be integrated to increase the accuracy of emotion identification.

Many Physiological signals are vital for emotion, according to several emotion theories. Ekman [5] proved that a particular physiological pattern might be linked to a specific emotion. For feature extraction, two types of data are employed in this study: EEG and peripheral. Frequency and amplitude are two parameters that may be used to define EEG. The EEG signal may be categorized into different bands based on frequency: Alpha, Beta, Gamma, Delta, and Theta. The frequency range of the above bands are 1 - 4 Hz., 4 - 8 Hz., 8 -12.5 Hz., 12.5 - 28 Hz. and 30-40 Hz, respectively.

9.2. MULTIMODAL EMOTION DATABASE

The availability of multimodal datasets is an essential factor in evaluating the performance of a pattern recognition system. The following are the multimodal databases that are currently available:

- HUMAINE database: This is an audiovisual database with gestures.
- Audiovisual database of Belfast naturalists.
- Smartkom combines audiovisual and gestural input.
- Salas audiovisual database.

Emotions Modelling in 3D Space

Abstract: In this study, we have discussed emotion representation in two and three-dimensional space. The three-dimensional space is based on the three emotion primitives, *i.e.*, valence, arousal, and dominance. The multimodal cues used in this study are EEG, Physiological signals, and video (under limitations). Due to the limited emotional content in videos from the DEAP database, we have considered only three classes of emotions, *i.e.*, happy, sad, and terrible. The wavelet transforms, a classical transform, were employed for multi-resolution analysis of signals to extract features. We have evaluated the proposed emotion model with standard multimodal datasets, DEAP. The experimental results show that SVM and MLP can predict emotions in single and multimodal cues.

Keywords: Arousal, DEAP database, Dominance, EEG, Multiresolution analysis, Support vector machine, Valence.

10.1. INTRODUCTION

Emotion modelling is an effective mechanism to implement emotions in various domains like Artificial Intelligence, Robotics, Psychology and Human-Computer Interaction. Modelling emotions enables machines to exhibit human being like behavior. With the advancement of technology, affect recognition is shifting from distinct emotions toward continuous two or three-dimensional space. Now, researchers focus on the investigation of the dimensional model of emotion. The Dimensional model deals with complex emotions that can be represented by two or three-dimension space. Generally, a 2D model has two dimensions: valence and arousal (in positive and negative axes). A 2D model is typically represented by four quadrants. Some scholars have used various terminologies for distinct emotion dimensions.

Whissell [1] developed a two-dimensional emotion model based on appraisal and activation. We looked at some studies on emotion representation in two-dimensional spaces. We reviewed the two-dimensional emotion model in this chapter. The limitations of the two-dimensional model and the needs of the three-dimensional emotion model are discussed. We will examine how emotions are represented in three-dimensional VAD spaces. After that, a 3D emotion graph is

presented by expressing many emotions in three-dimensional space. In separate portions of this chapter, emotion predictions from multimodal signals are also defined. Finally, we conclude this chapter by presenting the essential findings and contributions.

To represent emotions in three-dimensional space, we used the DEAP database [2], an extensive database encompassing more than twenty-five emotions [3, 4, 5] have all effectively employed the DEAP database in emotion recognition experiments.

The significant contributions of this work for predicting multimodal emotion in VAD space are as follows:

1. It presents the approach for affect prediction regarding valence, arousal, and dominance based on facial expression, EEG, and peripheral cues.
2. It also predicts the correlation between emotion dimensions and demonstrates significant performance improvement.
3. It compares state-of-the-art machine learning techniques *i.e.*, Support Vector Machine and Multilayer Perceptron for continuous affect prediction.

10.2. AFFECT REPRESENTATION IN 2D SPACE

Whissell [1] proposed a two-dimensional emotion model based on the appraisal and activation of two values. They depict the location of emotive words in a (-3; +3) Evaluation-Activation (EA) space. Each of the six Ekman's fundamental emotions is represented as a 2D EA space. Each emotion has specific affective weights to be determined by the emotion analyzer.

Fig. (10.1) illustrates the 2D valence-arousal space as proposed by Whissell. It is fascinating that the emotion of 'joy' exists in isolation, even though it should be associated with cheerfulness and pleasure. Furthermore, the emotions 'surprise' and 'joy' belong to the same quadrant 2. Furthermore, intermediary states between 'joy' and the rest of the emotions are rarely detected by the Whissell space. As a result, we may conclude that the existing theory cannot be validated using this two-dimensional model. Now the open question is whether two-dimension space suffices to represent all emotions or if there is a need for more dimensions to represent emotions efficiently.

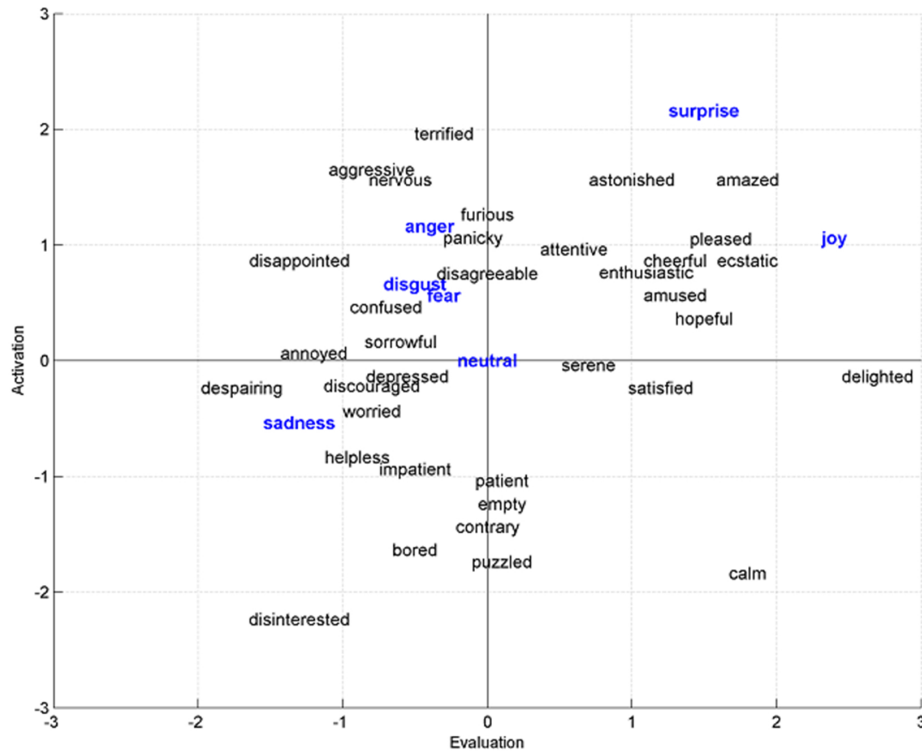


Fig (10.1). EV space as proposed by Whissell.

In an article [6], an anonymous researcher described a mood and emotion monitoring experiment. The resultant map of emotions shows that the two-dimensional mood measure does an excellent job of capturing emotional states. The positive and powerful feelings, such as pride, excitement, and joy, ended up in the space's top right quadrant. Irritation, anguish, and worry were in the third quadrant, which was the reverse of the first. The weak emotions are around the center of the arousal scale. However, stronger emotions are in the higher and lower regions of the plane, depending on valence.

As can be seen from Fig. (10.2), 'contentment,' 'affection,' and 'sadness,' are in the same valence axis in the VA space, which is not feasible according to existing theory. However, because the feeling 'affection' lies in the cheerful group, it must be close to the word 'joy.' This model also incorporates restricted emotions.

In valence-arousal space, this model is unable to reflect emotions appropriately. So, once again, the issue arises: do we need more dimensions to express emotions correctly?

SUBJECT INDEX

A

Adaptive 34, 21
 clustering technique 34
 mechanisms 21
 Affect-based stress reduction 9
 Affective 10, 17
 information complementary 17
 intelligence 10
 Affective computing 8, 9
 robots 8
 technologies 9
 Algorithms 9, 21, 25, 49, 76, 87, 104, 137
 multiple pre-processing 137
 Analogue to digital converter (ADC) 79
 Analog voltage signal 79
 Analysis 5, 17, 42, 63, 73, 111, 128, 140
 cross-modal factor 42, 111
 human motion 5
 machine emotion 17
 multi-resolution 73, 128
 Analysis of variance (ANOVA) 63, 140
 Anticipate complicated emotions 6
 Applications 70, 73, 122
 contemporary 122
 line detection 73
 real-world 70
 Arousal 22, 23, 24, 25, 33, 34, 35, 36, 42, 128,
 131, 132, 142, 143, 144
 and valence self-assessment 35
 music videos 35
 Array-based sensor's accuracy 78
 Artificial Intelligence 2, 5, 14, 22, 52, 67, 128,
 140
 Artificial Neural Networks 22, 52, 67, 140
 Asynchronous feature level fusion 111
 Audio-visual 36, 98, 107, 111
 cues fusion framework 107
 databases 98
 descriptions 36
 emotion recognition 111
 Auditory 21, 75, 98

Facial Expression 75
 Automatic propensity 18

B

Bayesian inference methods 52
 Bayesian Networks 3, 14, 43, 51
 dynamic 43, 51
 Biometrics 60, 71, 75
 iris recognition 60
 Biosignals 41, 66, 115
 peripheral 41, 66
 Blood Pressure 1, 10, 13
 Blood Volume Pressure (BVP) 6, 16, 53, 65,
 115, 136
 Brain activity 1, 116
 detecting continuous electrical 116
 Burmese-Mandarin corpus's database 100

C

Canon-bard theory 30
 Charge Coupled Device (CCD) 78
 electric 78
 Cinematographic principles 34
 CMOS sensors 78
 CNN-LSTM pooling 93
 Cognitive 2, 7, 15, 136
 appraisal theory 136
 processes 2, 15
 rehabilitation 7
 Cohn-Kanade database 43, 80, 85, 92, 93
 Communication 5, 6, 9, 13, 14, 15, 17, 18, 42
 emotional 9
 human-machine 6
 non-verbal 14, 42
 spoken 15
 verbal 18, 42
 Complimentary metal-oxide semiconductor
 (CMOS) 78
 Computers 1, 5, 8
 intelligent 8

Gyanendra K. Verma

All rights reserved-© 2023 Bentham Science Publishers

vision techniques 5
 Computing 5, 6, 8, 55, 61, 75, 132
 community 132
 devices 5
 emotional 6, 8
 power 75
 resources 55, 61
 Conditions 2, 7, 9, 10, 15, 16, 18, 60
 arousing 2, 15
 emotional 7, 9, 10, 16, 18
 eye 60
 Convolutional neural network 145
 Curvelet transform 22, 59, 71, 72, 73, 75, 84,
 139
 discrete 139
 technique 139

D

Darwin's theory 33, 37
 Data 2, 14, 22, 52, 60, 76, 121, 122
 emotional information 22
 multisensory 2, 14
 noisy 2, 76
 sensory 52
 training 121, 122
 transmission 60
 Data fusion 3, 14, 56, 61
 sensor-level 55, 61
 DBN speech recognition 53
 DEAP 23, 25, 53, 98, 115, 118, 123, 126, 128,
 129, 133, 136, 137, 140, 144, 145
 database 115, 118, 123, 128, 129, 133, 136,
 137, 139, 140, 144
 multimodal dataset 145
 Deep learning (DL) 22
 Devices, input 5
 Diagnostic technique 46
 Digital 78, 79
 imaging 78, 79
 photography 79
 Dimension(s) 33, 34, 36, 66, 67, 82, 128, 129,
 130, 133, 135, 145, 116
 geometric 116
 reduction methods 66
 Discrete 41, 53, 66, 97, 119, 123, 137, 139
 curvelet transform (DCT) 139
 wavelet transforms (DWT) 41, 53, 66, 97,
 119, 123, 137
 DS theory 53

image segmentation 53
 Dynamic bayesian network (DBN) 34, 43, 51,
 53

E

Early fusion methods 51
 EEG 23, 24, 25, 53, 65, 66, 115, 116, 117,
 118, 123, 125, 126, 128, 132, 137, 139,
 145
 emotion recognition 125
 measures voltage changes 115
 EEG and peripheral physiological 118, 123,
 132, 139
 data 132
 signals 118, 123, 139
 Effects 2, 6, 15, 19, 76
 automatic 19
 emotional 6
 environmental 76
 perceptual 2, 15
 Efficient modeling of emotions 30
 Ekman's emotion 35, 136
 Electroencephalogram 22
 Electrical 79, 116
 activity 116
 brain activity 116
 charges 79
 Electrocardiogram 41, 66, 115
 Electrode 41, 53, 66, 115, 116, 119
 placement approach 116
 dermal activity (EDA) 41, 53, 66, 115, 119
 Electroencephalograms 6, 16, 115, 118
 Electromyography measures 117
 Electrooculogram 16, 41, 66, 115
 EMG skin reaction 65
 Emotional 4, 9, 17, 64
 body language research 17
 integration 64
 intelligence 4
 strategies 9
 Emotional contact 9, 18
 human-human 18
 Emotion(s) 5, 6, 9, 13, 22, 42, 67, 111, 134
 expression 13
 features 67
 framework 5
 generation 6
 identification systems 22, 42, 111
 network 134

- perception, automatic 13
 - psychological 9
 - Emotion recognition 3, 5, 21, 22, 97
 - methods 21
 - research 3, 5
 - systems 22, 97
 - Environment 8, 9, 43, 52
 - virtual teaching 9
 - Environmental conditions 76
 - Equal error rate (EER) 44
- F**
- Face action coding system (FACS) 19, 104
 - Facial 4, 5, 7, 10, 17, 19, 20, 80, 104, 138
 - activity units 19
 - blood flow imaging 10
 - emotions 4, 5, 7, 17, 20, 80, 104, 138
 - Facial action 19, 20
 - coding method 20
 - unit (FAU) 19
 - Facial expression 5, 18, 76
 - analysis 5
 - fundamentals 18
 - images 76
 - recognition system 76
 - Fisher discriminate analysis (FDA) 66, 82, 97
 - Fourier transform 70, 71, 72
 - Frequency, vibration 101
 - Function 4, 63, 67, 68, 69, 71, 72, 83, 84, 87, 120, 122, 125
 - hyperbolic tangent 120
 - linear 68
 - mathematical transformation 71
 - non-linear activation 83, 122
 - sigmoid activation 122
 - Fusion 49, 50, 52, 53, 54, 55, 56, 59, 60, 61, 62, 65, 73, 97, 107, 108, 109, 110, 115, 140
 - adaptive 50
 - architectures 50, 140
 - audio-visual 52
 - framework 59, 65, 73, 115, 140
 - process 49, 50, 60
 - sensor-level 54, 55, 61
 - techniques 49, 50, 52, 56

G

- Gabor transform 22
- Galvanic skin response (GSR) 16, 41, 53, 66, 115, 119, 126, 136, 139
- Ganapathy 24, 145
- Gaussian 76, 77, 90, 94
 - process 94
 - white noise 90
- Gaussian noise 75, 87
 - conditions 87
- Genetic algorithm 22
- Ground truth data 139

H

- Half total error rate (HTER) 44, 45
- Heisenberg uncertainty principle 70
- Hidden Markov model (HMM) 3, 14, 25, 35, 42, 43, 93, 126
- Hierarchical dynamic bayesian network (HDBN) 34
- Human action 41, 66, 104
 - recognition community 41, 66
 - retrievals 104
- Human 1, 3, 14, 16, 17, 42, 53, 145
 - activity monitoring 53
 - centered approach 3, 14
 - computer interface 42
 - development research 17
 - machine interface 145
 - robot interaction (HRI) 1
 - social masking 16
- Human-computer interaction (HCI) 1, 3, 5, 13, 14, 16, 18, 98, 128
 - applications 16

I

- Illumination 21, 60
 - conditions 60
 - unregulated 21
- Image 21, 27, 72, 120, 141
 - backdrop complexity 21
 - processing techniques 72
- Image noise 78, 79
 - digital 79
- Information 40, 42, 54, 61
 - processing 40, 42

sources 54, 61
Information fusion 53, 54, 55
semantic 49, 53, 54, 55
Instance-based non-parametric learning 83
method 83
Intensity 34, 136
emotional 34
Interactions 3, 17, 42
man-machine 3
nonlinear 42
regular human-human 17
Interactive evolutionary algorithm 104

K

Kernel technique 70

L

Lagrange multipliers 69, 121
Lange's theory 30
Lazarus' cognitive-mediational theory 37
Lazy algorithm 122
Learning 52, 67, 145
algorithm 67
approaches, cutting-edge machine 145
Learning-based
fusion technique 52
method beats 52
Linear 41, 65, 100
discriminant analysis (LDA) 65
prediction cepstral coefficients (LPCC) 41, 65, 100
predictive coding (LPC) 100
Logarithmic relative power energy (LRPE) 118, 120
Log frequency power coefficient (LFPC) 41, 65, 100
Logistic regression algorithm 126

M

Machine learning 46, 120, 123, 141
algorithms 46
methods 123
techniques 120, 141
Machine recognition of facial expressions 21
Mandarin database 100
Mapping body gestures 17

Matthews correlation coefficient (MCC) 85
Mean 40, 45
absolute error (MAE) 45
square error (MSE) 40, 45
Multi-Resolution Analysis (MRA) 70, 78
approaches 70
based method 98
Multilayer perception (MLP) 82, 83, 84, 85, 102, 103, 105, 106, 108, 110, 121, 123, 143, 144, 145
Multimodal 42, 43, 49, 50, 59, 61, 63, 65, 67, 69, 71, 73, 97, 98, 117, 126, 137
databases 117, 126
data processing 50, 59, 137
emotion databases 98, 117
fusion 42, 49, 50, 59, 61, 63, 65, 67, 69, 71, 73, 97, 98
information fusion 43, 49, 50, 61
technique 42
Multimodal emotion 118, 145
forecasting 145
Multiresolution 41, 66, 76, 77
approaches 76, 77
techniques 41, 66
Multiresolution analysis (MRA) 59, 70, 71, 75, 76, 97, 98, 123, 128, 138, 140, 144, 145
technique 70
transform-based 76
Muscles 19, 117
facial 117
Muscular 1, 13, 104
activity 1, 13, 104
contraction 117
Music video clips 132

N

Natural language processing (NLP) 40
Neural network 42, 83, 94, 104, 111, 121
architecture 121
multi-layered feed-forward 83
Noise 3, 20, 60, 75, 76, 77, 78, 79, 80, 82
additive 77
channel 60
induced distortion 77
photon 79
pixel non-uniformity 79
thermal 79
Noisy database 76

Non-parametric learning technique 122
 Non-verbal communication techniques 2

O

Object-tracking method 61
 Oresteia database 118

P

Parrots' hypothesis 32
 Perceptual linear prediction (PLP) 41, 65, 100
 Physiological 6, 16, 17, 18, 34
 markers 6, 17
 muscle movements 34
 reactions 16, 18
 Pitch log energy 101
 Pixel response non-uniformity (PRNU) 79
 Power energy 115, 120, 126
 absolute 115
 features 126
 Preprocessing techniques 82
 Principal component analysis (PCA) 22, 65, 93, 94
 Probabilistic inference 3, 14, 51
 Processing 30, 82
 digital signal 82
 methods 30
 Properties 41, 42, 65, 72, 100, 101
 acoustic 100
 geometrical 72
 visual 42
 Psycho-evolutionary theory 32
 Psychological methods 9

R

Recognition systems 21, 30, 34, 35, 60
 adaptive expression 21
 multimodal emotion 34
 Regression 35, 70, 120
 combining support vector 35
 efficient non-linear 70
 Relative 45, 118, 120, 126
 absolute error 45
 power energy (RPE) 118, 120, 126
 Respiratory volume (RV) 6
 RML 97, 98, 99, 100, 101, 102, 104, 105, 106, 107, 108, 109

database 97, 98, 99, 100, 101, 102, 104, 105, 106, 107, 108, 109
 emotion database 100

Root mean squared error (RMSE) 45

S

SAM method 132
 Satellite imaging 75
 SAVEE database 100
 Scales 71, 72, 84, 131, 132, 139
 microscopic 71
 Schachter's theory 33
 Self-assessment manikins (SAM) 132
 Sensors 5, 9, 15, 17, 49, 50, 51, 54, 60, 61, 78, 79
 camera 60
 failures 51
 fusion 49
 Sensory modalities 3
 Sequential forward floating selection (SFFS) 139
 Signal 15, 30, 60, 61, 49, 60, 70, 71, 72, 77, 79, 80, 123, 128, 138, 141
 electrical 60
 emotional 77
 input 30, 61
 paralinguistic 15
 Skin conductance (SKC) 6, 41, 66, 115, 117, 119
 measurements 41, 66, 115, 119
 response (SCR) 41, 66, 115
 Social 8, 7
 connections 8
 neurosciences 7
 Spectral clustering method 34
 Strategy 43, 44, 79, 144
 asynchronous feature-level fusion 43
 Support vector machine (SVM) 53, 82, 93, 94, 97, 98, 102, 103, 105, 106, 108, 110, 111, 120, 123, 126, 128, 140, 141, 145
 Regression (SVR) 42, 67, 68
 SVM technique 120
 Systems 1, 2, 3, 6, 14, 16, 33, 42, 43, 44, 45, 60, 61, 76, 87, 93, 104, 109, 110, 111
 autonomic nervous 16
 combined 45, 109
 cue-based 109
 facial action coding 104
 human neurological 16

neurophysiological 33

T

Technologies 1, 8, 14, 128, 141

discrete wavelet transform (DWT) 141

Thermal agitation 79

V

VAD space 67, 129, 133, 134, 135, 145

Valence, arousal, and dominance (VAD) 2, 5,
34, 36, 128, 131, 133, 136, 137, 139,
140, 144, 145, 146

Video 52, 99, 136

data 52, 136

emotion database 99

Virtual interfaces 7

Visual information 3, 41, 66, 78, 139

W

Wavelet 70, 81, 84, 85, 90, 92, 97, 111, 119,
141, 142, 145

analysis 81, 90, 92, 97, 111, 119, 142

curvelet analyses 84

Wavelet coefficients 141

function 85

transformations 70, 145



Gyanendra K. Verma

Gyanendra K. Verma is currently an Assistant Professor (senior grade) in the Department of Information Technology, National Institute of Technology Raipur, India. He completed his B. Tech. in 2006 from Harcourt Butler Technical University (formerly HBTI) Kanpur, India and M. Tech. & Ph.D. from the Indian Institute of Information Technology (IIITA) Allahabad, India, in 2009 and 2016, respectively. He has teaching and research experience of more than 10 years in the area of Computer Science and Information Technology, with a special interest in Image Processing, Speech and Language Processing and Human-Computer Interaction. He has collaborated actively with researchers in several other disciplines of computer science, particularly on Medical Imaging problems at the SILP laboratory, IIIT Allahabad. He has served on roughly fifty conference and workshop program committees and served as the Organizing Chair for MIND 2019. His research on the application of Wavelet Transform in Medical Imaging and Computer Vision problems have been cited extensively. He is a senior member of Institute of Electrical and Electronics Engineers (IEEE) and Association for Computing Machinery (ACM); reputed International Technical Societies.