# ASSESSING GAMEPLAY EMOTIONS FROM PHYSIOLOGICAL SIGNALS
## A FUZZY DECISION TREES BASED MODEL

**Joseph Onderi Orero** $^{*a}$**, Florent Levillain**$^{b}$**, Marc Damez-Fontaine**$^{a}$**, Maria Rifqi**$^{a}$
**and Bernadette Bouchon-Meunier**$^{a}$

$^a$*Laboratoire d'Informatique de Paris 6 (LIP6), Université Pierre et Marie Curie, France*
$^b$*Laboratoire Cognitions Humaine et Artificielle (CHART), Université Paris 8, France*

## ABSTRACT

As video games become a widespread form of entertainment, there is need to develop new evaluative methodologies for acknowledging the various aspects of the player's subjective experience, and especially the emotional aspect. Video game developers could benefit from being aware of how the player reacts emotionally to specific game parameters. In this study, we addressed the possibility to record physiological measures on players involved in an action game, with the main objective of developing adequate models to describe emotional states. Our goal was to estimate the emotional state of the player from physiological signals so as to relate these variations of the autonomic nervous system to the specific game narratives. To achieve this, we developed a fuzzy set theory based model to recognize various episodes of the game from the user's physiological signals. We used fuzzy decision trees to generate the rules that map these signals to game episodes characterized by a variation of challenge at stake. A specific advantage to our approach is that we automatically recognize game episodes from physiological signals with explicitly defined rules relating the signals to episodes in a continuous scale. We compare our results with the actual game statistics information associated with the game episodes.

**Keywords:** Emotion Recognition, Video Games, Physiological Signals, Fuzzy Sets.

## 1.  INTRODUCTION

In the recent years video games have enjoyed a dramatic increase in popularity, the growing market being echoed by a genuine interest in the academic field, with an attempt to justify the existence of video game theory in academia [1]. In this context of flourishing technological and

---

$^*$**Corresponding author:**104, Avenue du Président Kennedy, 75016 Paris, France. joseph.orero@lip6.fr

theoretical efforts, it is regrettable that effective methods of evaluation of the gaming technology still lag behind. Until now, the evaluation of video games has mainly relied on traditional methods through subjective self-reports such as questionnaires [2], interviews and focus groups [3]. These subjective reports have limitations, one of them being that they only generate data at the moment the question is asked, and not through a continuous process. Also, subjects are known to be poor in self reporting their behaviors during game situations [4]. Since the main goal of video games is to entertain through a continuous renewal of the user's interest, controlling the emergence of certain affective states is a key determinant in achieving a truly immersive experience. Therefore, in evaluating video games, it would be more appropriate to assess dynamically the emotional responses in relation to variations in challenge difficulty and to specific tasks at hand in the game.

Secondly, it has a particular interest regarding the field of developing affective human-computer interaction applications. The video games industry appears as a laboratory where we can develop innovative devices of interaction that could be used as a standard for designing affective-oriented interfaces keeping track of users's emotional states.

It is thus important to develop algorithms to automatically compute the affective state during the play process. In this study, we addressed the possibility of recording physiological measures on players involved in an action game, with the main objective of developing adequate models to describe emotional states. Specifically, our main goal was to estimate the emotional state of the player from physiological signals so as to: (i) relate these variations of the autonomic nervous system to the specific game narratives and (ii) differentiate these game narratives according to their level of challenge from physiological signature independent of the player.

The rest of the paper is organized as follows: In Section 2 we give an overview of affective computing. In Section 3, we outline the methods of classification and physiological measures used. We then describe the experimental setup in Section 4 and give our results in Section 5. Finally we give conclusions and future work in Section 6.


## 2. AFFECTIVE COMPUTING

Ever since Picard's highlight [5] on the subject, affective computing has received considerable interest from many researchers in Human Computer Interaction (HCI). In particular the recognitions of emotions from physiological signals has become a major research focus in the recent past [6]. Among a vast range of possible ways to access a user's emotional responses such as facial gesture or voice recognition through video and audio recording, subjective self report measures, physiological measures stand out. They grant an access to non conscious and non reportable processes [7] and may to a certain extent be unobtrusively monitored [8]. In addition, physiological measures may provide a fine-grained account of the rise, surge or extinction of an emotional response.

Considerable research progress has been made in the use of physiological measures to reliably predict specific emotions [9, 10, 11]. However, recognition of emotions through physiological signals is not yet a clear systematic task. There is lack of standardized methodology in the analysis

of physiological data for affective system in terms of models, measures, features and algorithms to use [12]. The difficulty to uniquely map physiological patterns onto specific emotion types is due to the nature of these bio data. Physiological data tends to vary considerably from one person to another and may even display considerable differences within individuals on different occasions [9, 10, 11]. The aim of this work was to address some of these difficulties.

First, one of the main difficulty in emotion recognition is to ensure that people are really in the expected emotions. Researchers have employed various methods of eliciting emotions such as guided imagery technique, movie clips, math problems and music songs [9, 13, 11]. However, use of video games to elicit emotions is more advantageous in many ways. Particulary, due to their interactive nature, video games tend to promote a natural sense of immersion, the player being concerned by the consequences of its own actions, especially when mastering a skill is at stake. In this work, we took advantage of this propensity to immerse the user with the hope that we could drive specific physiological reactions corresponding to variations in the challenge proposed by the game at different points.

Secondly, emotions are better represented in fuzzy terms rather than discrete categories. As Goulev [14] already pointed out, even human ability to identify emotions from physical appearance is rarely hundred percent and as such, it is necessary to express in fuzzy terms mapping emotions from physiological data. Moreover, the physiological data from sensors is itself imperfect such that it is difficult to express the results in crisp terms [15]. Also, one could express more than one emotional state category at a particular time. A fuzzy set theory based approach could better represent both the multi-state complexity and uncertainties as well as imperfections in emotion recognition.

## 3. CLASSIFICATION AND PHYSIOLOGICAL MEASURES

### 3.1. Classification Methods

In emotion recognition from physiological signals, various methods of classification have been used in the past such as linear discriminant analysis, k-nearest-neighbor (KNN), multilayer perceptron network and decision trees (DT) [9, 13, 10, 16, 11]. A comparison of various methods using optimal features [10] seems to suggest that results depend on the method chosen and the nature of the experiment. In particular, decision trees have been found to perform comparably well [16]. A major superiority characteristic of decision trees is that they give explicitly defined rules used in classification. This is important in our case to explicitly know the relationship between the physiological signals attributes and the emotional states.

Secondly, as we have already noted, it is preferable to use a fuzzy sets theory based approach in classification. In fuzzy set theory [17], a fuzzy set is represented by a membership function, $\mu_A : A \rightarrow [0, 1]$, indicating the degree to which an element belongs to a given set $A$. This is a contrast to $\{0, 1\}$ in a crisp set, in which an element can only belong to a given set (membership value of 1) or not (membership value of 0). It is interesting for this kind of recognition to express the emotional states in a continuous scale of values $[0, 1]$. Therefore, we choose to use fuzzy

decision trees (FDT) in this work. Fuzzy approach has been used in the past to model emotions from physiological signals [14, 18]. Particularly, Mandryk and Atkins [18], developed a fuzzy logic model of rules that were used to transform arousal and valence based on psychophysiology literature. The aim of our experiment was to automatically generated these rules that define the psychophysiology relations from the fuzzy decision trees based on the players' physiological data. We used Marsala's Salammbô Fuzzy Decision Tree [19] and compared our results with Quilan's C4.5 decision trees [20] and KNN [21] .

## 3.2. Physiological Signals and Features

Based on previous literature, we chose to collect galvanic skin response (GSR) which is a measure of the conductivity of the skin. GSR is considered as an effective correlate to arousal [22, 23, 24] and has been extensively used in the domain of affective computing [8, 25, 18]. We also collected heart rate (HR) through a measure of cardiovascular activity. HR may differentiate between positive and negative emotions [26, 27, 23] although the possible correlation between heart rate and valence remains debated [28, 29].

After collecting the signals, the raw signals were expressed in values $[0, 1]$ for each subject's data. Then as in [9], we calculated from each of these signals, six statistical features (the mean, the standard deviation, and the mean of the first and the second absolute difference of the raw signal and the normalized signal) for each game episode segment. However, in our case, for each feature, we subtracted the subject's baseline value. The baseline value was calculated by taking the value of each player's signal from the two minutes period preceding the beginning of the game. This was to account for the variations from subject to subject in signal values and also to inform us whether the attribute value increases or decreases on playing the particular game episode.

## 3.3. Feature Selection

A key research question is to determine which set of measures and features of these measures provide the optimal combination for discriminating the different emotional categories [5]. There is no agreement on a particular feature selection method [12], it highly depends on the nature of the data and method used in classification. In this work we tested Sequential Backward Search (SBS), Sequential Forward Search (SFS) and Sequential Floating Forward Search (SFFS) which have been found to perform well in this kind of data [9, 10, 30, 11]. However, there was no difference in results between them, the optimal feature set was the same for all the three. This could be expected since, unlike other classifiers, decision trees also perform feature selection by selecting the best attribute to discriminate the classes in each node of the tree. Nevertheless, to be able to rank features in order of their relevance in classification, we utilized SFS.

## 4. THE EXPERIMENTAL SETUP

### 4.1. Game Episodes

In this study we tested a game belonging to a popular genre in the game industry, Halo3, which is a First-Person Shooter (FPS). This game appears as one of the most immersive genre in that it propels the player at the heart of action through a first-person perspective. To segment the game

session in different episodes we used specific events scripted inside the game. Scripted events are triggered by certain actions from the player, they are mandatory and are typically used by game designers either to mark transitions between different periods of the game, to inform the player about a certain state of the world or to draw her attention toward key information. They manifest through dialogs, visual and/or sound effects (e.g. the rumor of a battleship afar indicating the proximity of the enemy). By taking advantage of these scripted events we could delimitate three different combat episodes. Each episode is characterized by distinct geographical cues and corresponds approximatively to a game level in the traditional acceptation where boundaries of a level are defined by the completion of a goal (in our case the complete "cleaning" of a zone).

Before playing the game, participants were told to rest for two minutes during which a movie was presented, depicting abstract forms evolving smoothly and slowly so as to help the participant to relax. This period was considered as a baseline for the subsequent physiological recording. The first combat sequence (Combat1) is characterized by a succession of short skirmishes. The overall difficulty in this episode is low as it is almost possible to progress without shooting while fellow soldiers take care of the enemies. This is followed by the second combat session (Combat2) which is globally a bit more difficult than Combat1. The last combat (Combat3) begins when the player arrives at an enemy site. This is the most intense episode with enemies outnumbering and taking cover. In order to contrast combat episodes with more relaxing periods, we distinguished two additional sessions consisting in transitory periods (REST), triggered by the death of the last enemy in a level, wherein the player is able to rest, collect ammunitions and explore the environment.

### 4.2. Participants and Setting

The experiment was conducted at the LUTIN (Laboratoire des Usages en Technologies d' Information Numérique) in the Cité des Sciences et de l'Industrie, Paris, France. Twelve male participants aged between 18 to 40 were recruited from visitors of the Cité des Sciences et de l'Industrie to participate in the experiment. No specific expertise in the field of video games was required, although we selected participants able to manipulate a gamepad and to orient themselves in a virtual environment. Halo3 was played on a Microsoft Xbox 360 on a 32-inch LCD television. A camera captured the TV screen. Participants were seated at approximatively one meter from the screen, they were explicitly told not to move and keep the game pad onto their laps in order to avoid any muscular artefact in the physiological recording. All participants played to the game Halo3, at the easiest to difficult level. Participants were allowed to play at their own pace, until the experimenter told them to quit the game. The game ended when participants reached a specific point in the game (end of Combat3).

To collect the physiological measures we used the Biopac MP35 acquisition unit and the software BSLPro to visualize the data. To record the HR, we measured the electrocardiography (ECG) through a Einthoven derivation II placing pre-gelled surface electrodes on the ankles and on the wrist. We recorded GSR using surface electrodes attached with Velcro Straps that were placed on two fingers of the left hand. The fingers wearing the electrodes remained wedged under the gamepad. In order to synchronize the video with the physiological data, we used sound markers emitted at the beginning and at the end of the baseline video which were sent to the acquisition

unit. ECG and GSR data were collected at 200 Hz. As noisy ECG data may produce failures in computing the HR, we inspected the HR data and corrected manually every erroneous samples.
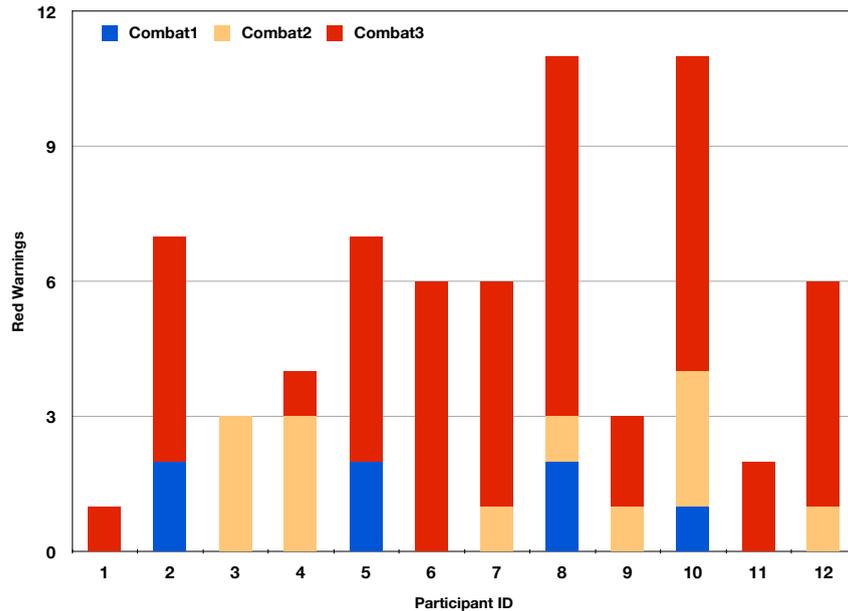
## 5. RESULTS



**Figure 1**: Numbers of Red Warnings during Combat1, Combat2 and Combat3 for each participant

As we have already noted, our main goal in this experiment was to automatically discriminate episodes of the game characterized by high challenge from those characterized by low challenge based on physiological signals. The three different combat episodes could be characterized in term of their level of difficulty, based on certain information provided by the player activity. Specifically, we took advantage of a warning signal triggered each time the player is in danger, i.e. when his energy bar is fully depleted. This kind of information is a fair indicator of the level of intensity of an episode, as it reflects the amount of stress placed upon the player in a circumstance of likely death. As can be seen in Figure 1, there is a global and significant increase of the number of warning signals in Combat3 when compared to Combat1. However, there is only a slight increase of warning signals in Combat2 when compared to Combat1. As a consequence, we chose to contrast between Combat1 and Combat3.

As earlier explained, we calculated 6 features from each physiological signals and utilized these attribute values as input to the classifiers. Although the output from the fuzzy decision tree was in continuous values $[0, 1]$, we deffuzified the values into $\{0, 1\}$ for quantitative analysis and to compare the results with other classifiers. In the first attempt, a classifier was trained by leaving out one subject's data (test sample) against the rest of the subjects' data (training sample) using the aggregate of the entire episode (3 to 5 minutes). This was to give us a global recognition model of various game episodes. But in a related way, we considered a second alternative, in which we subdivided each game episode into 20 segments (approximately 10 seconds length). This was to

validate if the model could apply on a continuous scale at various times during the game episodes. Half (six) of the subjects' samples was used for training and the other half for testing. In the analysis, we give the results of both alternatives.

First, we performed machine learning to differentiate between game combats of high challenge from those of low challenge with reasonable success as shown in Table 1 and Table 2. In discriminating these episodes, the GSR (mean amplitude and the mean of the second difference of the normalized) signal was the most relevant when compared across the classifiers using both the aggregate of the entire episode and the 10 seconds segments samples.

**Table 1**: Low/High Challenge Episodes Pre-labeled Vs Predicted Results (Aggregate of Entire Episode)

| - | Low | High | Salammbô FDT | Quilan C4.5 DT | KNN(k=1) |
|---|---|---|---|---|---|
| Low | 8 | 4 | 66.67% | 50.00% | 66.67% |
| High | 1 | 11 | 91.67% | 83.33% | 58.33% |
| Average | | | 79.17% | 66.67% | 62.50% |

**Table 2**: Low/High Challenge Episodes Pre-labeled Vs Predicted Results (10 Seconds Segments)

| - | Low | High | Salammbô FDT | Quilan C4.5 DT | KNN(k=10) |
|---|---|---|---|---|---|
| Low | 202 | 38 | 84.17% | 85.00 % | 80.00% |
| High | 49 | 191 | 79.58% | 77.08% | 86.67% |
| Average | | | 81.88% | 81.04% | 83.33% |

Secondly, we tried to discriminate between episodes wherein the player is involved in combat and transitory periods (REST). As already noted, the transitory periods (REST) came after the death of the last enemy in the combat episode and thus index the moments the players feels relieved to have accomplished a goal. The player is not involved in any shooting while moving forward to start the next combat episode. During our machine learning process, we managed to differentiate these episodes (combat and REST) with 68.75% and 70.4% success as shown in Table 3 and Table 4 respectively. It turned out that, in addition to GSR (standard deviation) signal, the HR (the mean of the first and second difference) signal was very relevant in discriminating the REST/combat categories across the classifiers using both the aggregate of the entire episode and the 10 seconds segments samples.

**Table 3**: Combat/REST Episodes Pre-labeled Vs Predicted Results (Aggregate of Entire Episode)

| - | Combat | REST | Salammbô FDT | Quilan C4.5 DT | KNN(k=1) |
|---|---|---|---|---|---|
| Combat | 19 | 5 | 79.17% | 41.70% | 70.08 % |
| REST | 10 | 14 | 58.33% | 54.20% | 70.08 % |
| Average | | | 68.75% | 47.95% | 70.08 % |

Generally, segmentation of episodes seems to have increased the classification rate across all the classifiers. Although using the aggregate of the entire episode gives us a global classification model, it drastically reduces the sample size for good machine learning. The performance of classifiers is also comparable especially when episodes are segmented, if judged by classification

**Table 4**: Combat/REST Episodes Pre-labeled Vs Predicted Results (10 Seconds Segments)

| - | Combat | REST | Salammbô FDT | Quilan C4.5 DT | KNN(k=10) |
|---|---|---|---|---|---|
| Combat | 377 | 103 | 78.54% | 76.0% | 77.3 % |
| REST | 181 | 299 | 62.29% | 66.5% | 69.8 % |
| Average | | | 70.42% | 71.25% | 73.55 % |

into discrete categories. However, fuzzy decision trees will be more applicable as they can be used on a continuous scale.

## 6. CONCLUSIONS

In this study we tested a model to automatically recognize specific episodes in an action game. Our aim was to relate a set of physiological signals features to specific game episodes. By taking advantage of the game narratives and the corresponding variations in the challenge proposed by the game through its different levels, we trained fuzzy decision trees to predict the episodes and generate rules that map the features to these periods of the game.

Although we couldn't differentiate clearly between combat episodes and transitory episodes where the player is not involved in a fight, we managed to identify with considerable success different combat episodes characterized by a variation in the challenge at stake. Our results point to galvanic skin response as a very relevant measure of the level of arousal of a player since it best discriminated between very intense episodes and not so intense episodes. Similarly, we found that, in addition to galvanic skin response, heart rate was also very relevant measure in discriminating between shooting and the subsequent rest episodes. We anticipate that variations of these signals could index the moments the players feels relieved when a goal has been accomplished.

However, several issues inhibited us from obtaining optimal results in this experiment. One issue regards to the players level of expertise. As we couldn't control properly the level of expertise of our participants, we may have failed to obtain homogeneous samples. Confronted to the same level of difficulty, players with sensibly different playing background may react very differently. For instance a well trained player would react positively to an especially intense episode, as it may represent a fair challenge to him. Inversely, a beginner would consider that the same episode is outrageously punishing and thus, frustrating. Therefore, we would hope to obtain an even better level of recognition by testing people with the same amount of experience, and possibly the same kind of reaction confronted to an increase in intensity.

Also, the indexation method we used to differentiate the game episodes may have failed to be the most accurate in differentiating clearly between combat episodes and transitory episodes where the player is not involved in a fight. It is possible that psychological boundaries of an episode (i.e. when a period of the game is subjectively felt as starting and ending) only partially coincide with objective markers such as scripted events. In consequence, it would be useful to consider the possibility to index events according to psychological markers.

Much is still to be done before getting access to the structure of the player's emotional pro-

cesses. In particular, we need to improve our methods in order to define episodes with respect to psychological states and their variations. As a next step toward recognizing levels of affective involvement in a game, we would first need to consider more physiological measures. The present study did not exhaust all the possible signals available and improvements can be expected from the consideration of a large amount of physiological features. Additionally, we would need to differentiate physiological signatures not only from a level of intensity we assume they elicit, but from the players' direct assessment. We could therefore evaluate episodes based on various subjective scales (e.g. how frustrated I am, how concentrated I am, how fun is my experience, etc.) so as to define the episodes on several dimensions. Altogether, the road map for the forthcoming investigation of affective states in video game will get through a clear definition of the most relevant dimensions to account for the emotional response we target, as well as a thorough examination of the best physiological features. This way, we hope for a truly systematic affective recognition procedure to be incorporated to the games evaluation routines.

## REFERENCES

[1] Wolf, M and Perron, B. Introduction. In *The Video Game Theory Reader*. Routledge, 2003.

[2] Nel, F, Damez, M, Labroche, N, and Lesot, M.-J. Automated video games evaluation based on the template formalism. In *International Conference on Information Processing and Management of Uncertainty in Knownledge-Based Systems, IPMU*, 2008.

[3] Grassioulet, Y. A cognitive ergonomics approach to the process of game design and development. Master's thesis, Educational Technologies (STAF), TECFA, Psychology And Educational Sciences Department, University of Geneva, 2002.

[4] Pagulayan, R. J, Keeker, K, Wixon, D, Romero, R, and Fuller, T. User-centered design in games. In *Handbook for human-computer interaction in interactive systems*. Lawrence Erlbaum Associates, Inc. Mahwah, NJ, 2002.

[5] Picard, R. W. *Affective Computing*. The MIT Press, Cambridge, Massachusetts, London, England, 1997.

[6] Fairclough, S. Fundamentals of physiological computing. *Interacting With Computers*, 21:133–145, 2009.

[7] Cacioppo, J. T, Tassinary, L. G, and Bernston, G. Psychophysiological science: Interdisciplinary approaches to classic questions about the mind. In *Handbook of Psychophysiology*. Cambridge University Press, 2007.

[8] Scheirer, J, Fernandez, R, Klein, J, and Picard, R. Frustrating the user on purpose: A step toward building an affective computer. *Interacting With Computers*, 14:93–118, 2002.

[9] Picard, R, Vyzas, E, and Healey, J. Toward machine emotional intelligence: Analysis of affective physiological state. *IEEE Transactions Pattern Analysis and Machine Intelligence*, 23:1175–1191, 2001.

[10] Wagner, J, Kim, J, and André, E. From physiological signals to emotions: Implementing and comparing selected methods for feature extraction and classification. In *IEEE International Conference in Multimedia and Expo*, 2005.

[11] Kim, J and André, E. Emotion recognition based on physiological changes in music listening. *IEEE Transactions on Pattern Analysis And Machine Intelligence*, 30:2067–2083, 2008.

[12] Arroyo-Palacios, J and Romano, D. M. Towards a standardization in the use of physiological signals for affective recognition systems. In *Measuring Behavior*, 2008.

[13] Nasoz, F, Alvarez, K, Lisetti, C, and Finkelstein, C. Emotion recognition from physiological signals for presence technologies. *International Journal Of Cognition, Technology And Work*, 6, 2003.

[14] Goulev, P. *An Investigation into the Use of AffectiveWare in Interactive Computer Applications*. PhD thesis, Imperial College London, 2005.

[15] Bouchon-Meunier, B. *Aggregation and Fusion of Imperfect Information*. Physica-Verlag, Spring-Verlag Company, 1998.

[16] Rani, P, Sarkar, N, and Adams, J. Anxiety-based affective communication for implicit human machine interaction. *Advanced Engineering Informatics*, 21:323–334, 2007.

[17] Zadeh, L. Fuzzy sets. *Information Control*, 8:338–358, 1965.

[18] Mandryk, R. L and Atkins, M. A fuzzy physiological approach for continuously modeling emotion during interaction with play technologies. *International Journal of Human-Computer Studies*, 65:329–347, 2007.

[19] Marsala, C. Fuzzy decision trees to help flexible querying. *KYBERNETICA*, 36:689–705, 2000.

[20] Quinlan, R. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers, San Mateo, CA, 1993.

[21] Cover, T and Hart, P. Nearest neighbor pattern classficiation. *IEEE Transactions Information Theory*, 13:21–27, 1967.

[22] Lang, P. The emotion probe studies of motivation and attention. *American Psychologist*, 50(5):372–385, 1995.

[23] Bradley, M, Greenwald, M. K, and Hamm, A. Affective picture processing. In *The Structure of Emotion*. Hogrefe and Huber Publishers, Toronto, 1993.

[24] Detenber, B, Simons, R, and Bennett, G. Roll' em! the effects of picture motion on emotional responses. *Journal of Broadcasting and Electronic Media*, 21:112–126, 1998.

[25] Ward, R and Marsden, P. Physiological responses to different web page designs. *International Journal of Human-Computer Studies*, 59:199–212, 2003.

[26] Winton, W, Putnam, L, and Krauss, R. Facial and autonomic manifestations of the dimensional structure of emotion. *Journal of Experimental Social Psychology*, 20:195–216, 1984.

[27] Papillo, J. F and Shapiro, D. Principles of psychophysiology: Physical, social, and inferential elements. In Tassinary, L, editor, *The Cardiovascular System*. Cambridge University Press, 1990.

[28] Neumann, S and Waldstein, S. Similar patterns of cardiovascular response during emotional activation as a function of affective valence and arousal and gender. *Journal of Psychosomatic Research*, 50:245–253, 2001.

[29] Ritz, T and Thöns, M. Airway response of healthy individuals to affective picture series. *International Journal of Psychophysiology*, 46:67–75, 2002.

[30] Yang, R and Liu, G. Emotion feature selection from physiological signal based on bpso. In *ISKE-2007 Advances in Intelligent Systems Research*, 2007.