# A Study on Automotive Human Vehicle Interaction Using Gesture Recognition Technology

**Santosh Naik[1], Abhishek.H.R[2], Ashwal.K.N[3], Balasubramanya.S.P[4]**
Department of Information Science and Engineering, Jain University, Bangalore, India
[1]santoshrcnaik@gmail.com
[2]abhihrsirsi@gmail.com
[3]ashwal.navada@yahoo.com
[4]balu_samaga@yahoo.com

**Abstract:** In order to establish a more natural and effective communication with virtual systems, hand gesture recognition appears as a suitable means. The main objective of gesture recognition is to create a system which can recognize specific human gestures and use them to convey proper message for device control or system and perform the task accordingly. The approaches on gesture recognition can be mainly classified as Data Glove Based approaches and Vision Based approaches. This paper describes Human and Machine Interface using gesture recognition as an alternative to the user using Vision Based technology. Few of the gesture recognition technologies are discussed in brief along with its application. However the main objective of the author's interest in gesture recognition for automotive human vehicle interaction is to evaluate possible application by using hand gestures in cars (our focus is on GPS). This paper concludes with our approach on gesture recognition.

**Keywords:** Virtual System, Gesture Recognition, Data Glove Based approach, Vision Based approach, Human and Vehicle Interface

## INTRODUCTION

Gesture recognition is basically recognizing the gestures of a human, involving the hands, arms, face, head, and/or body. It is of utmost importance in designing an intelligent, efficient and accurate human–computer interface.

"Have you seen the remote?" "I left it on the table after watching my matinee show". "It is not here, I will miss the news again because of you!" In the near future, such heated discussions over remote control won't disturb the harmony of the house. Not because they will place it correctly but because soon remote controls will be the objects of the past. Technology has finally reached that dimension when our hands will take over the job and replace them by directly communicating with the computer or television. For instance, in order to delete a folder or file from the computer, place your palm on it, and throw it like a paper in a dustbin. Even while using the microwave oven to bake a cake, waving our hands in the air like a magician would serve as a command for the oven. While some of us might be thinking of it being a futuristic vision, some of us have already experienced it through what we call **"gesture recognition technology"**. Since the time that the computer revolution started, human computer interaction has always been attempted to improve. Computers have now become an integral part of our lives and hence their usage should be as trouble-free as talking to someone is. Earlier the way humans interacted with this smart machine was either through keyboard or a mouse. But now attempts are being made to make the man-machine interaction as natural as possible. Fulfilling this requirement is the popular touch screen technology which is soon expected to be replaced by the **gesture recognition technology** [1]. However the major concern of this technology is user acceptance.

## OVERVIEW

In addition to the theoretical aspects, any practical implementation of gesture recognition typically requires the use of different imaging and tracking devices or gadgets. These include instrumented gloves, body suits, and marker based optical tracking. Traditional 2-D keyboard, pen, and mouse-oriented graphical user interfaces are often not suitable for working in virtual environments. Rather, devices that sense body position and orientation, direction of gaze, speech and sound, facial expression, and other aspects of human behaviour or state can be used to model communication between a human and the environment. For a real-time application, the expectation is to obtain the best possible images of the hand gesture within the lowest possible time. Some experiments have been conducted with the purpose of defining the best configuration for imaging the hand. This configuration includes, among others, the relative position of the hand and the camera, the influence of the integration time of the camera, the amplitude threshold, and the lighting conditions of the environment, the surrounding objects and the skin colour. Some of these parameters are discussed in [2] in detail while some of them are discussed below in brief.

**Sensitivity to integration time**

The integration time is the length of time that the pixels are allowed to collect light. Not setting the appropriate integration time results in a loss of information (saturated pixels) while imaging the hand. Pixels receiving more light than expected are saturated and do not provide any information in the output image. Ideally, the more the light is collected without saturation, the better the hand gesture will be imaged. Saturation occurs due to excessive signal and/or background light. Thus one should look for the highest possible integration time without having any saturated pixels in the hand blob. The appropriate integration time is a function of distance and has been determined empirically during an experiment where hand images are collected at different distances with different integration times.

**Sensitivity to lighting conditions**

We can make use of SR4000 in our approach. The SR4000 is an active sensor that emits infrared light and records the reflected radiation through a filter. In order to check whether any additional desk lamp could disturb our applications, an experiment has been conducted by [2] where a spot light has been placed in front of the SR4000 camera. It has been noticed that part of the light coming from the lamp passes through the filter and is recorded by the camera. This creates some saturated pixels but it doesn't prevent imaging the hand if the number of saturated pixels is insignificant with respect to the size of the hand segment. It has also been noticed that the light at the ceiling as well no light in the room have any specific influence on the hand images. As a concluding remark about the lighting conditions for imaging the hand by making use of a range camera, avoiding any additional light directly placed in front of the camera is clearly advisable.

**Sensitivity to surrounding objects**

The SR4000 produces hand images which are highly dependent on reflexive surrounding objects. These objects may cause multiple reflections and/or multiple paths of some light rays and consequently generate some hanging pixels in the images which do not represent any physical object in the reality. In case some of these pixels appear between the hand and the camera, they should be discarded while extracting the hand information.

**HAND DETECTION AND GESTURE RECOGNITION**

Gesture recognition is the process by which gestures made by the user are made known to the system. It can also be explained as the mathematical interpretation of a human motion by a computing device. After detecting the hand, the task is to find which gesture, if any, is performed. Recognition of gestures consists of two steps: 1) capturing the motion and configuration/pose of fingers, hands, and/or arms, depending on the level of detail required (hereafter: hand), and 2) classify the captured data as belonging to one of the predefined gesture classes. A number of different devices have been applied in order to capture the data, e.g., magnetic devices,

accelerometers, and bend sensors, but in general the capturing is either performed by a glove-based system or by an optical-based system. Due to the non-intrusive nature of the latter, it is the focus in much gesture research. How this is done is described in the next section. The approach of gesture recognition is divided into two steps corresponding to two different algorithms, one which detects the number of outstretched fingers and one which handles the point and click gestures.

**Count the number of fingers**

From the observations made by [3], it follows that a very simple approach to counting the number of outstretched fingers is to do a polar transformation around the centre of the hand and count the number of fingers (rectangles) present in each radius. As the gestures are only performed while the hand is pointing upwards, only the interval [180±; 360±] is investigated. In order to speed up the algorithm the segmented image is sampled along concentric circles instead of doing polar transformation. The algorithm does not contain any information regarding the relative distances between two fingers. The reason for this is firstly because it makes the system more general, and secondly because different users tend to have different preferences depending on individual kinematics limits of their hands and fingers. The designed algorithm is thus robust to how the different gestures are performed, e.g., the three-finger gesture can also be performed by an outstretched ring finger, middle finger, and index finger. In fact, each gesture can be performed in a number of different ways.

Number of different configurations to perform a gesture $i$.

| Gesture: | $i$ | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|---|
| Configurations: | ($5Ci$) | 1 | 5 | 10 | 10 | 5 | 1 |

**Recognise the point and click gestures**

When the system recognises that only one finger is present using the algorithm, this is interpreted as a pointing gesture. The tip of the pointing finger is defined to be the actual position where the user is pointing at. This point is found as follows: for the consecutive radii used in [3] to classify the pointing gesture, the centre of the finger is found for each radius and these values are fitted to a straight line. This line is then searched until the final point is reached, i.e., the finger tip. Similar to a computer mouse we also need a click interaction which can be associated with the current position of the pointing gestures. For example, a click gesture can indicate that the virtual object currently being pointed at should be selected.

**HAND GESTURE RECOGNITION TECHNOLOGIES**

Reference [4] talks about gesture detection, tracking, recognition in detail. The discussion related to the functionalities and the corresponding algorithms are beyond the scope of this paper. Reference [5] describes centroid

7

features, Hidden Markov model (HMM), various databases in depth. Reference [1] describes various types of gesture recognition technologies in use currently. They are:

### Contact type

It involves touch based gestures using a touch pad or a touch screen. Touch pad or touch screen based gesture recognition is achieved by sensing physical contact on a conventional touch pad or touch screen. Touch pads and touch screens are primarily used for controlling cursors on a PC or mobile phones and are gaining user acceptance for point of sale terminals, PDA's, various industrial and automotive applications as well. They are already being used for automotive applications, and PDA's. User acceptance of touch-based gesture automotive systems technologies are relatively easier for the public to accept because they preserve a physical user interface.

### Device gesture technologies

Device-based techniques use a glove, stylus, or other position tracker, whose movements send signals that the system uses to identify the gesture. There is a body of research detailing gesture technologies [[6], [7]]. Reference [6] describes an example of an impractical technology for HVI: "hand gesture only interfaces with syntax of many gestures typically require precise hand pose tracking. A common technique is to instrument the hand with a glove which is equipped with a number of sensors which provide information about hand position, orientation, and flex of fingers" such as the data glove [8]. Raw data collection using instrumented gloves and trackers requires users to physically attach computer input devices to their hands.

The instrumented gloves report data values for the movement of the fingers; the amount of reported data values depends on the type of glove worn. The trackers are attached to the back of the hand or the upper wrist, depending on the type of glove worn, and give back data on the position and orientation of the hand in 3D space. First commercial hand tracker, Data Glove, used thin fibre optic cables running down the back of each hand, each with a small crack in it. Light is shone down the cable so when the fingers are bent light leaks out through the cracks. Measuring light loss gives an accurate reading of hand poses. Similar technique is used for wearable suits used in virtual environment applications. Though gloves provide accurate measurements of hand shape, they are cumbersome to wear, they are connected through wires, which restricts freedom of movement and they are expensive too. Various other kinds of systems are reported in literature for intrusive hand gesture recognition. Some uses bend sensor on the index finger, an acceleration sensor on the hand, a micro switch for activation. To reduce physical restriction due to the cables, an alternate technique used is to wear an ultrasonic emitter on the index finger and the receiver capable of tracking the position of the emitter is mounted on a Head Mounted Device (HMD). Wearing a glove is clearly not a practical proposition for automotive applications and wearable suits as used in virtual environment applications would be impractical also.

Due to space limitations, the reader is encouraged to refer to the existing literature on HMM evaluation, estimation, and decoding [[9], [10], [11], [12]]. A tutorial relating HMM's to sign language recognition is provided in the first author's master's thesis [13].

### Vision-based technologies

One of the main difficulties in using glove-based input devices to collect raw posture and gesture recognition data is the fact the gloves must be worn by the user and attached to the computer. In many cases, users do not want to wear tracking devices and computer-bound gloves since they can restrict freedom of movement and take considerably longer to set up than traditional interaction methods. As a result, there has been quite a bit of research into using computer vision to track human movement and extract raw data for posture and gesture recognition. A vision-based solution to collecting data for hand posture and gesture recognition requires four equally important components. The first is the placement and number of cameras used. Placing the camera(s) is critical because the visibility of the hand or hands being tracked must be maximized for robust recognition. Visibility is important because of the many occlusion problems present in vision-based tracking. The second component in a vision-based solution for hand posture and gesture recognition is to make the hands more visible to the camera for simpler extraction of hand data. The third component of a vision-based solution for hand gesture and posture recognition is the extraction of features from the stream or streams of raw image data; the fourth component is to apply recognition algorithms to these extracted features.

Vision-based gesture interpretation system as follows:-
(1) Video input from camera.
(2) Analysis of the video input.
(3) Recognition of the gesture.

Once the raw data has been collected from a vision-based data collection system, it must be analysed to determine if any postures or gestures have been recognized. Various algorithmic techniques for recognizing hand postures and gestures are discussed in [[14], [15], [16]].

The great advantage of vision based gesture recognition is the liberty that is given to the person in which gesture recognition is performed. As opposite to the device based gesture recognition that use gloves or accelerometers boxes, the vision based gesture recognition is a non-intrusive method. However, it can nevertheless be observed several disadvantages. Often, more computation power is required for vision based gesture recognition compared to device based gesture recognition. Another problem is the difficulty to perform vision-based gesture recognition when we have different backgrounds or different lighting conditions. There are two approaches to vision based gesture recognition:

8

*1) Model based techniques:* They try to create a three dimensional model of the users hand and use this for recognition. Some systems track gesture movements through a set of critical positions. When a gesture moves through the same critical positions as does a stored gesture, the system recognizes it. Other systems track the body part being moved, compute the nature of the motion, and then determine the gesture. The systems generally do this by applying statistical modelling to a set of movements.

*2) Image based methods:* Image-based techniques detect a gesture by capturing pictures of a user's motions during the course of a gesture. The system sends these images to computer-vision software, which tracks them and identifies the gesture. These methods typically extract flesh tones from the background images to find hands and then try and extract features such as fingertips, hand edges, or gross hand geometry for use in gesture recognition.

### E. Electrical field sensing:

Proximity of a human body or body part can be measured by sensing electric fields; the term used to refer to a family of non- contact measurements of the human body that may be made with slowly varying electric fields. These measurements can be used to measure the distance of a human hand or other body part from an object. This facilitates a vast range of applications for a wide range of industries.

## OUR APPROACH

We believe that a good gesture interface as a minimum requires a pointing and a click gesture. The primary argument for this is that many applications may be controlled through these two gestures. Furthermore, it should be easy for the user to remember how to perform these gestures. Thinking in terms of the pointing gesture, the most natural way to performing this is by an outstretched index finger. It is important that the system should recognise the gestures instantly also the system should be able to distinguish between purposeful gestures and accidental ones, unless proven otherwise we'd still be afraid that if we sneezed the system would do something screwy like automatically switch over to the heavy metal channel on satellite radio at full volume and blow out our eardrums. We design some predefined gestures to minimize the search space. To keep track of the finger trajectories, we design an efficient algorithm. Reference [5] focuses on various databases. In the LTI-gesture database, very good results can be obtained using a downscaled image of each video frame and tangent distance as a model of image variability. However HMM is most fundamentally and widely used for gesture recognition and mean shift algorithm is widely used for tracking. A detailed approach on HMM is discussed in [5]. Fig.1 shows the representation of our approach.



**Fig.1:** General Architecture of gesture recognition

Reference [17] describes the general approach.

### Data acquisition:

The data acquisition component is responsible for processing the received data and then transmits them to the gesture manager. First, a set of filter is used to optimize the data. For example, the position/orientation information is very noisy due to dependence of lighting conditions. Thus, orientation data that exceed a given limit are discarded as improbable and replaced with their previous values.

### Gesture manager

The gesture manager is the principal part of the recognition system. This library maintains a list of known postures. The system tries to match incoming data with existing posture. This is done by first looking for the best matching fingers constellation. Five dimensional vectors represent the bend values of the fingers and for each posture definition the distance to the current data is calculated. Then, the position/orientation data is compared in a likewise manner. Finally, in this gesture recognition system, a gesture is just a sequence of successive postures. For example, let's consider the detection of a "click" gesture. This gesture is defined as a pointing posture with outstretched index finger and thumb and the other fingers flexed, then a tapping posture with half-bent index finger.
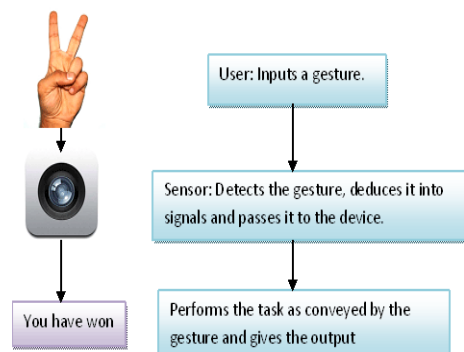


**Fig.2:** Typical stages of gesture recognition

9

**Fig.3:** Gesture recognition using counting the number of fingers

In our approach, it may be noted that the gestures are irrespective of the hand you use; either left or right. Fig.2 shows the typical process of gesture recognition. Gesture could include operation of the interior light; as the hand is detected approaching the light switch the light could switch on. If the hand is detected approaching the light switch again it would switch off, thus the hand movement to and from the device being controlled could be used as a pre-emptive gesture. It could even adjust the brightness of the light by using basic hand gestures. Swirling your fingers in the clockwise direction would increase the brightness and similarly would decrease the brightness of the light if done so in anti-clockwise direction. Similar observations can be made in stereo and radio. Flicking your fingers on one direction would skip to the next song in the stereo or would switch over to the next channel in the radio. If same is done in the opposite direction, it would revert back to the previous track in the stereo or would switch to the previous channel in the radio. Controlling of volume in both the cases can be done by swirling the fingers just like adjusting the brightness in the former. Possible context sensitive gestures to indicate yes/no or accept/reject could be thumbs-up and a thumbs-down. The same yes/no gestures could also be used to accept or reject prompts for automatic navigation re-routing. For example, if the advanced navigation system has been informed of an accident ahead it could ask the driver if he/she wishes to be automatically re-routed to avoid possible delays. Similarly, if low fuel is detected, the system could ask the driver if he/she would like to be automatically routed to the nearest fuel station. It may also be possible to simply use one of the yes/no gestures for other functions. The gesture system consists of a visual prompt for displaying the menus and an SR4000 sensor to recognise the gesture. If a person fanned his hand in front of his face, the gesture system could detect this and interpret that he is too hot and would like to cool down. Visual prompt could be offered to the person to ask if he/she would like to adjust the temperature by displaying a set of temperatures on the visual prompt. The person can use any of the techniques to select the desired temperature. Firstly the technique which makes use of counting the number of fingers and perform the task accordingly based on its order. Fig.3 shows the possible gestures which can be recognised by the system. For example if the person fanned his hand in front of his face, the system would display the set of temperatures. Here the camera detects the number of outstretched fingers. If the system detects any of the gestures mentioned in the Fig.3, the corresponding numbered temperature present in the list would be selected and the AC would be switched on.  The second one being using the point and click gestures. Here the SR4000 sensor

would track the trajectories of the fingers and move the pointer accordingly. To select the required temperature, a click gesture can be made use. We also make use of GPS which would probably be more helpful rather than taking out i-phones and surfing for a nearest restaurant or a café. Natural gestures for i'm hungry or thirsty initiate a dialogue of "do you want to be directed to the nearest café or restaurant?" in the visual prompt. The GPS could be used to display the nearest café or restaurant. If the person wishes to be directed to that location, a thumbs-up gesture would direct him the location. If he wants to check for the restaurant or a café which is a bit distant than this, then a thumbs-down gesture would discard the route to the nearest restaurant or a café and would automatically route him to the next restaurant or a café which is a bit distant than this. The applications are only limited by our ability to find a natural gesture to initiate the required meaningful dialogue. Likewise, an eye contact would probably carry out certain tasks such as opening or closing the window.

**APPLICATIONS**

In general, gesture recognition has many potential uses. Hand gesture recognition can be used in scientific research, construction, and any situation where user and machine must be separated, such as dangerous areas suited to mechanical exploration. Gesture recognition can be used in creating adaptable assistive technologies for persons with various physical and sensory handicaps. There are several communicative applications, such as sign language recognition. Pure facial gesture recognition can be used to great effect for purposes related to personal security and law enforcement. Another application for hand postures and gestures is control of audio and video devices. Freeman and Weissman have developed a system to control a television set by hand gestures [18]. Using an open hand, the user can turn the television on and off, change the channel, increase and decrease the volume, and mute the sound. Other applications that could use hand gestures are control of a VCR, a stereo, or a whole room [19]. After considerable analysis of all the possible selective themes and functions, it was found that each set of gestures could fit into the following application domain classifications [1]:-
(1) Pre-emptive gestures
(2) Function associated gestures
(3) Context sensitive gestures
(4) Global shortcut gestures
(5) Natural dialogue gestures
Each classification is now discussed.

**Pre-emptive gestures**
A pre-emptive natural hand gesture occurs when the hand is moving towards a specific control type or device and the detection of the hand approaching is used to pre-empt the drivers intent to operate a particular control.

10

### Function associated gestures

Function associated gestures are those gestures that use the natural action of the arm/hand to associate or provide a cognitive link to the function being controlled.

### Context sensitive gesture

Context sensitive gestures are natural hand gestures that are used to respond to driver prompts or automatic events. Possible context sensitive gestures to indicate yes/no or accept/reject could be thumbs-up and a thumbs-down.

### Global shortcut gestures

Global shortcut gestures are in fact natural symbolic gestures that can be used at any time, the term natural refers to the use of natural hand gestures that are typically used in human to human communications. It is expected that hand gestures will be selected whereby the user can easily link the gesture to the function being controlled.

### Natural dialogue gestures

Natural dialogue hand gestures utilise natural gestures as used in human to human communication to initiate a gesture dialogue with the vehicle, typically this would involve two gestures being used although only one gesture at any given time.

Gesture recognition has wide-ranging applications [20] such as the following:

(1) Developing aids for the hearing impaired;
(2) Enabling very young children to interact with computers;
(3) Designing techniques for forensic identification;
(4) Recognizing sign language;
(5) Medically monitoring patients' emotional states or stress levels;
(6) Lie detection;
(7) Navigating and/or manipulating in virtual environments;
(8) Communicating in video conferencing;
(9) Distance learning/tele-teaching assistance;
(10) Monitoring automobile drivers' alertness/drowsiness levels, etc.

### CONCLUSION

As we have seen, gesture is a movement of a limb or a body as an expression of thought or feeling. And by using gesture recognition system we can make interface with computer using gesture of human body, typically hand movements. In human recognition technology, a camera reads the movements of the human body and communicates the data to a computer that uses the gestures as input to control devices or applications. Data gloves are used to sense the gesture. There are two basic methods of gesture recognition.

(1) Constellation method
(2) Hidden Markov model

The constellation model is somewhat simpler than HMM because it includes actual samples of gestures from already stored sample gestures and it contains less mathematical approach. But disadvantage of this method is it is very slow and takes around two minutes to recognize one gesture. On other hand HMM is a rich statistical model and most widely used. This model can able to give high probability equations by which we can able to represent hidden states in a given sequence. Also this model is a double stochastic process. Except for the long wait, we give it thumbs up.

### REFERENCES

[1][Online] Available: http://www.engineersgarage.com/articles/ gesture-recognition-technology

[2] Hervé Lahamy and Derek Litchi, "*Real-Time Hand Gesture Recognition using Range Cameras*", University of Calgary, 2500 University Dr NW, Calgary, Alberta, T2N1N4, Canada.

[3]Moritz Störring, Thomas B. Moeslund, Yong Liu, and Erik Granum, "*COMPUTER VISION-BASED GESTURE RECOGNITION FOR AN AUGMENTED REALITY INTERFACE*", IN PROC 4th IASTED International Conference on VISUALIZATION, IMAGING, AND IMAGE PROCESSING, Marbella, Spain, Sep 2004, pages 766-771.

[4]X. Zabulisy, H. Baltzakisy, A. Argyroszy, "*Vision-based Hand Gesture Recognition for Human-Computer Interaction*", Heraklion, Crete, Greece.

[5]Prof. Dr.-Ing. H. Ney, "*Appearance-Based Gesture Recognition*", Dipl.-Inform. D. Keysers, January 2005.

[6]Anderson, James A, "*An Introduction to Neural Networks*", Bradford Books, Boston, 1995.

[7] Butterworth B and G. Beattie, "*Gesture and Silence as Indicators of Planning in Speech. In Recent Advances in the Psychology of Language*", Campbell and Smith (eds.), Plenum Press, New York, 1978.

[8] Huang, X. D., Y. Ariki, and M. A. Jack, "*Hidden Markov Models for Speech Recognition*", Edinburgh University Press, Edinburgh, 1990.

[9] L. Baum, "*An inequality and associated maximization technique in statistical estimation of probabilistic functions of Markov processes. Inequalities*", 3:1-8, 1972.

[10] X.D. Huang, Y. Ariki, and M. A. Jack, "*Hidden Markov Models for Speech Recognition*", Edinburgh University Press, 1990.

[11] L. R. Rabiner and B. H. Juang, "*An introduction to hidden Markov models*", IEEE ASSP Magazine, pages 4-16, January 1986.

[12] S. Young, "*HTK: Hidden Markov Model Toolkit*" V1.5. Cambridge Univ. Eng. Dept. Speech Group and Entropic Research Lab. Inc., Washington DC, 1993.

[13] T. Starner, "*Visual recognition of American Sign Language using hidden Markov models*", Master's thesis, MIT Media Laboratory, February 1995.

[14] Kadous and Waleed, "*GRASP: Recognition of Australian Sign Language Using Instrumented Gloves*", Bachelor's thesis, University of New South Wales, 1995.

[15] Starner and Thad, "*Visual Recognition of American Sign Language Using Hidden Markov Models*", Master's thesis, Massachusetts Institute of Technology, 1995.

[16] Watson and Richard, "*A Survey of Gesture Recognition Techniques*", Technical Report TCD-CS-93-11, Department of Computer Science, Trinity College Dublin, 1993.

[17] Damien Zufferey, *"Device based gesture recognition"*, Department of Informatics (DIUF), University of Fribourg, Bd de Perolles 90, CH-1700 Fribourg, Switzerland.

[18] Freeman, William.T and Craig D. Weissman, *"Television Control by Hand Gestures"*, Technical Report, Mitsubishi Electronic Research Laboratories, TR-94-24, 1994.

[19] Kohler and Marcus, *"Special Topics of Gesture Recognition Applied to Intelligent Home Environments"*, In Proceedings of the International Gesture Workshop'97, Berlin, 285-297, 1997.

[20]      C. L. Lisetti and D. J. Schiano, *"Automatic classification of single facial images"*, Pragmatics Cogn., vol. 8, pp. 185–235, 2000.