

Hand Gesture to Speech Translation for Assisting Deaf and Dumb

Rajatha Prabhu
Department of ECE
SJBIT
Bangalore, India

Harshitha B
Department of ECE
SJBIT
Bangalore, India

Madhushree B
Department of ECE
SJBIT
Bangalore, India

Dr. K. R. Nataraj
Guide and HOD of ECE
SJBIT
Bangalore, India

Abstract-Communication is the integral part of life. About 360 million people in the world are suffering from hearing impairment and 32 million of these are children, and their life is not as easy as it is for human without barrier. This project presents the Sign Language Recognition system capable of recognizing hand gestures by using MATLAB. The proposed technique has 4 modules such as: pre-processing, segmentation, feature extraction, gesture recognition and gesture to voice conversion. Different features are extracted such as Eigenvalues and Eigen vectors which are used in recognition. The Principle Component Analysis (PCA) algorithm is applied for gesture recognition and a recognized gesture is converted into text and voice format. The proposed technique helps to minimize communication barrier between deaf/mute and normal people.

I. INTRODUCTION

Sign language is a language which mainly uses manual communication to convey meaning, as opposed to acoustically conveyed sound patterns. This can involve simultaneously combining shapes of hands, orientation and movement of the hands, arms or body & facial expressions to express a speaker's thoughts. In order to facilitate communication between hearing impaired & hearing people, sign language interpreters are usually used. Such activities involve considerable effort on the part of the interpreter, as sign languages are distinct natural languages with their own syntax, varied from any spoken language.

MATLAB (matrix laboratory) is a multi-paradigm numerical computing environment and fourth-generation programming language. A proprietary programming language developed by MathWorks, MATLAB allows matrix manipulations, plotting of functions and data, implementation of algorithms, creation of graphical user interfaces (GUI), and interfacing with programs written in other languages. Using MATLAB we can process video with functions and system objects that read and write video files, perform feature extraction, motion estimation and object tracking and display video.

II. RELATED WORKS

Nasser H.D et.al [1] considers this approach where the key features extracted are SIFT (Scale Invariant Feature Transform) key-points. They further constructed a grammar from a sequence of hand postures for detecting dynamic gestures.

In [2] a basis for usage of Hidden Markov Models (HMM) is established by drawing an analogous relationship between speech recognition and gesture recognition. HMMs can be used to model time series data, and here the movement of the hand along the coordinate axis is tracked and each direction is taken as a state. This paper makes use of a lexicon of forty gestures and achieves an accuracy of 95 percent. It also states its disadvantage that as the lexicon grows the need to describe the hand configuration along with hand trajectory also will grow making the designing of HMM more complex and time consuming. We needed a method to describe dynamic gesture in a simpler way.

In [3] the system uses an intrinsic mobile camera for gesture recognition and acquisition; gesture acquired is processed with the help of Algorithms like HSV model (Skin Color Detection), Large Blob Detection, Flood Fill and Contour Extraction. The system is able to recognize one handed sign representation of the standard alphabets (A-Z) & numeric values (0-9). The output of this system is very efficient, consistent and of high approximation of gesture processing and speech analysis.

The paper [4] focuses on vision based hand gesture recognition system by proposing a scheme using a database driven hand gesture recognition based upon skin color model approach and thresholding approach along with an effective template matching using PCA. Initially, hand region is segmented by applying skin color model in YCbCr color space. In the next stage, thresholding is applied to separate foreground and background. Finally, template based matching technique is developed using Principal Component Analysis (PCA) for recognition.

In [5] Human computer interaction (HCI) & sign language recognition (SLR), aimed at creating a virtual reality, 3D gaming environment, helping the deaf-mute people etc., extensively exploit the use of hand gestures. Segmentation of the hand part from the other body parts and background is the primary need of any hand gesture based application system; but gesture recognition systems are usually plagued by different segmentation problems, and by the ones like coarticulation, recognition of similar gestures.

The primary aim of the work [6] is to design & implement a low cost wired interactive glove, interfaced with a computer running MATLAB or Octave,

with a high degree of accuracy for gesture recognition. The glove maps the orientation of the hand and fingers with the help of bend sensors, Hall Effect sensors and accelerometer. The data is then transmitted to computer using automatic repeat request as an error controlling scheme.

The algorithm devised in [7] is capable of extracting signs from video sequences under minimally cluttered & dynamic background using skin color segmentation. It distinguishes between static and dynamic gestures & extracts the suitable feature vector which are classified using Support Vector Machines (SVM). Speech recognition is built upon standard module -Sphinx.

[8] This paper presents the Sign Language Recognition system capable of recognizing 26 gestures from the Indian Sign Language (ISL) by using MATLAB. The proposed system having 4 modules such as: pre-processing and hand segmentation, feature extraction, sign recognition and sign to text and voice conversion. Segmentation is done by using image processing. Different features are extracted such as Eigen values and Eigen vectors which are used in recognition. The Principle Component Analysis (PCA) algorithm was used for gesture recognition & recognized gesture is converted into text and voice format.

This paper [9] presents an algorithm of Hand Gesture Recognition by using Dynamic Time Warping methodology. The system consists of three modules: real time detection of face region and two hands regions, track the hands trajectory both in terms of direction among consecutive frames as well as distance from the centre of the frame and gesture recognition based on analyzing variations in the hand locations along with the centre of the face. The proposed technique overcomes not only the limitations of a glove based approach but also most of the vision based approach concern illumination condition, background complexity and distance from camera which is up to two meters by using Dynamic Time Warping which finds the optimal alignment between the stored database & query features, improvement in recognition accuracy is observed compared to conventional methods. In [10] a Wireless data glove which is a normal cloth driving glove fitted with flex sensors is used along the length of each finger and the thumb. Mute people can use the gloves to perform hand gesture and it will be converted into speech so that normal people can understand their expression. A sign language usually provides sign for whole words. It can also provide sign for letters to perform words that don't have a corresponding sign in that sign language. In this paper, Flex Sensor plays the major role, Flex sensors are the sensors whose resistance changes depending on the amount of flexion. Here the device recognizes the sign language Alphabets and Numbers. It is in the process of developing a prototype to reduce the communication gap between differentiable and normal people. The program is in embedded C coding. Arduino software is used to observe the working of the program in the hardware circuitry which is designed using microcontroller and sensors.

III. DESIGN AND IMPLEMENTATION

In the project we generate the extract the skin color information from the video frames and use Principle Component Analysis (PCA) with the Euclidean distance as the classifier for the classification by using the Eigen values and vectors of the query image and the database images. Here in the proposed methodology uses MATLAB to process the frames of the query video to detect the gesture and by using indexing the particular hand gesture is recognized and the audio synthesizer would give the audio output for the detected action. Audio is pre-recorded for a particular hand gesture. The figure 1, shows the overview of the proposed system methodology for hand gesture recognition for providing artificial voice for the hearing impaired and mute people.

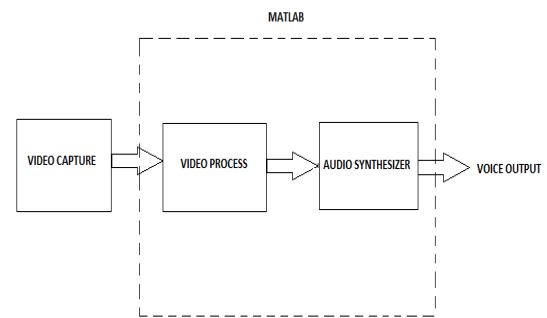


Figure 1 Proposed Methodology

A. Resize the images

The first step in the processing is to convert the query video frames into images and to resize it to a dimension, say, 280X280. Resizing of the image is necessary as the captured image would be of larger size and requires more memory to store the frames and to process it.

B. Extraction of the skin pixels

The resized images are then considered to extract the skin pixels from the image. The Grey world algorithm is used before the extraction of the skin pixels is done.

1) Grey-world algorithm

Color constancy is a technique for detection of color that does not depend on source of light. The source of light may add color casts in acquired images. To solve this problem a technique is to appraise the color of the predominant light and then, in the next stage, remove it. Once the color of light in individual channels is obtained each color pixel is normalized by a scaling factor.

One of the most commonly used simple methodology for estimating the color of light is Grey-World. This method provided a good result in practice if the average scene color is grey.

2) Grey-world assumptions

The Grey World Assumption is a white balance method that assumes that the scene, on average, is a neutral grey. Grey-world assumption holds good if the scene has good distribution of colors. Assuming this condition, the average reflected color is considered to be the color of the light. Hence, we estimate the illumination color cast by considering

the arithmetic mean color and compare it to grey. Grey world algorithm provides one with an estimate of illumination by computing the average of each channel of the image. One of the methods of normalization is that the mean of the three components is used as illumination estimate of the image. To normalize the image of a channel i , the pixel value is scaled by using equation (1),

$$S_i = \frac{avg}{avg_i} \dots\dots\dots(1)$$

where avg_i the channel mean and avg is the illumination estimate.

Another method is by normalizing to the maximum channel by scaling by

$$r_i = \frac{\max(avg_R, avg_G, avg_B)}{avg_i}$$

Another method is by normalizing to the maximum channel by using the equation(2) given below where m_i is calculated as,

$$m_i = \sqrt{(avg_r * avg_r + avg_g * avg_g + avg_b * avg_b)}$$

$$r_i = \frac{\max(m_R, m_G, m_B)}{m_i} \dots\dots\dots(2)$$

3) RGB to YCbCr conversion

After the application of grey-world algorithm the RGB images are converted to YCbCr as medical investigation has proven that the human eye has variable sensitivity to brightness and color. Hence the transformation of Red, Green and Blue color to YCbCr color space. The figure 2 shows the YCbCr color space.

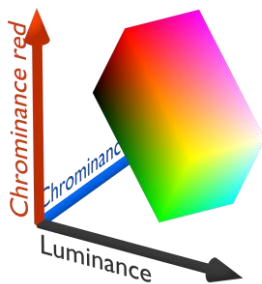


Figure 2 YCbCr color space

Y signifies the Luminance component, similarly Cb signifies the Chrominance-Blue along with the Cr component signifying Chrominance-Red. Grayscale form is analogous to the Y component of the actual image. Cb value is high in portions of the image having the blue color, the Cb and Cr values are low in portions of image with green and Cr value is high in portions of instances having shades of red. Medical research on the concept of the eye has resulted with the count of rods as 120 million, which are highly sensitive when compared to the cones whose count is 6-7 million. The rods are insensitive to color, whereas the cones provide the eyes

the sensitivity to color and are found to be situated close to the middle region.

The formula for transform RGB to YCbCr are given below,

$$Y = 0.3007R + 0.58593G + 0.11328B$$

$$Cb = 128 - 0.17187R - 0.33984G + 0.51171B$$

$$Cr = 128 + 0.51171R - 0.4296G - 0.08203B$$

The value range of Cb and Cr for skin colored pixel is given by,

$$Cb \geq 77 \ \& \ Cb \leq 127 \ \& \ Cr \geq 133 \ \& \ Cr \leq 173$$

The given skin color tone information includes a wide range of skin colored pixels.

Once the skin pixels are identified by using the range specified above, those pixels are marked with white pixels with a intensity of 255. The figure 3 shows the skin pixel extraction obtained.



(a) (b)

Figure 3 Extraction of the skin pixels in the image

C. Recognition of the gesture

The first step in the recognition of the gesture is to specify the path of the database images and the query images obtained from the frames of the video.

1) Principal Component Analysis (PCA)

Principal Component Analysis is a statistical method for performing the investigation of the correlations between a set of variables to find the fundamental architecture of those variables. It is also known as factor analysis. It is a analysis with nonparametric components and the output is particular and does not depend on any of the hypothesis information distribution.

The tasks that PCA can perform is forecasting, redundancy elimination, data compression, feature extraction, etc. Because

PCA is a classical methodology which works well in the linear hostname, applications having signal processing, image processing, system control theory, communications, etc. as linear models are compatible, Principal component analysis decreases the dimensionality of digital image but also keeps the image information and provides a compact features or compact representation of a digital image. The aim of the PCA technique is to transform the gesture images into a set of characteristics feature images called eigengesture. In recognition a query image is casted onto the lower-dimension gesture space traversed by the eigengetsure and then classified either by using a classifier or statistical theorem.

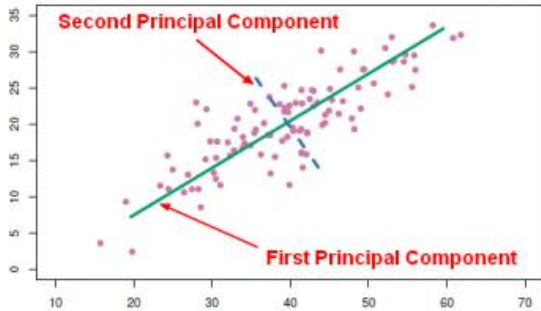


Figure 4 Example For PCA

The 2 principal components are defined as follows:

First principle component – the direction which maximizes variability of the data when projected on that axis.

Second principle component – the direction, among those orthogonal to the first, maximizing variability.

The principle components are eigenvectors of $A^T A$ and the eigenvalues are the variances.

2) *Eigen vectors and eigen values*

Eigenvalues and eigenvectors have being used widely in the matrices application in engineering and science. Image processing, vibration analysis, Control theory, electric circuits, and quantum mechanics are the areas of application. Many of the applications involving the use of eigenvalue sand eigenvectors is the process of transforming a given matrix into a diagonal matrix.

When we obtain a set of data points we can break the set into eigenvectors and eigenvalues. Eigenvectors and values exist in pairs: all of the eigenvectors have a corresponding eigenvalue. An eigenvector is a direction such as 45 degrees, vertical, horizontal, etc. An eigenvalue is a number, which specifies how much variance is there in the data in that direction, the eigenvalue is a number specifying how the data is distributed on the line. The principle component considered is the eigenvector with highest eigenvalue.

3) *Euclidean distance*

The Euclidean distance is the distance between the 2 pixels that is a straight-line. The figure 5 illustrates the Euclidean distance metrics.

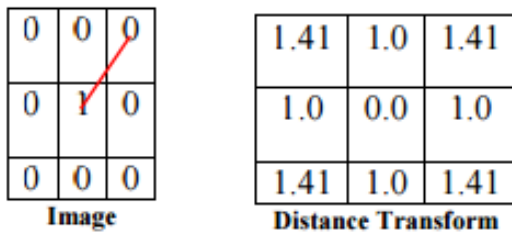


Figure 5 Illustration of the Euclidean distance metrics

Consider the two points P and Q in Euclidean spaces, two dimensional and P with the coordinates (p_1, p_2) , Q with the coordinates (q_1, q_2) . The hypotenuse of right angled triangle is the line with the endpoints as P and Q. The distance between 2 points p & q is given by the square root of the sum of the square of the differences between the corresponding coordinates of the points.

By Euclidean geometry in 2 dimensional space, the Euclidean distance between the 2 points $a = (ax, ay)$ & $b = (bx, by)$ is given by equation (3):

$$d(a, b) = \sqrt{(bx - ax)^2 + (by - ay)^2} \dots\dots\dots(3)$$

4) *Euclidean distance algorithm*

The least distance between a procured of column vectors in the code book matrix and a column vector x is computed by Euclidean distance algorithm. The algorithm computes the least distance to x and finds the column vector in code book that is nearest to x.

$$d(a, b) = |p - q|$$

$$\sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

$$= \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

In one dimensional space, the distance between 2 points, x_1 and x_2 , on a line is the absolute value of the difference between the 2 points by equation (4):

$$\sqrt{(X_2 - X_1)^2} = |X_2 - X_1| \dots\dots\dots(4)$$

In two dimensional space, the distance between $P = (p_1, p_2)$ and $q = (q_1, q_2)$ by equation (5):

$$\sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2} \dots\dots\dots(5)$$

5) *Euclidean function*

The input source data is a property class which will be converted within to a raster prior to the application of the Euclidean analysis. The figure 6 shows the illustration of Euclidean function.

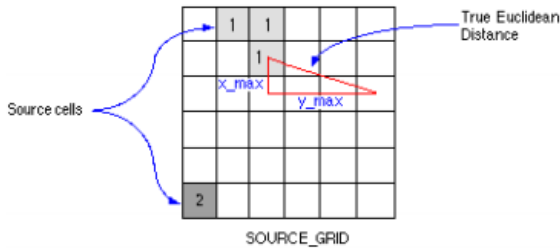


Figure 6 Euclidean Function

Euclidean distance for Images:

M X N images is easily analyzed in MN dimensional Euclidean space, called image space. The base e_1, e_2, \dots, e_{MN} are adopted to make up a coordinate system for the image space, where e_{kN+1} is for an ideal point source with unit intensity at $(k, 1)$. Hence image $x = (x_1, x_2, \dots, x_{MN})$, where x_{kN+1} is the grey level at the $(k, 1)^{th}$ pixel, is indicated as a point in the image space. An image whose grey levels are zero at all the points is the origin of the image space.

The metric coefficients g_{ij} $i, j = 1, 2, \dots, MN$, are given as

$$g_{ij} = \langle e_i, e_j \rangle = \sqrt{\langle e_i, e_i \rangle} \sqrt{\langle e_j, e_j \rangle} \cdot \cos \theta_{ij}$$

where the pointed brackets is the indication of scalar product, and θ_{ij} is the angle between e_i and e_j . Note that, if $\langle e_i, e_i \rangle = \langle e_j, e_j \rangle = \dots$, that is all the base vectors have the same length, then g_{ij} depends totally on the angle θ_{ij} . Given the metric coefficients, the Euclidean distance of two images x, y is given by

$$d_E^2(x, y) = \sum_{i,j=1}^{MN} g_{ij} (x^i - y^i)(x^j - y^j) = (x - y)^T G (x - y)$$

where the symmetric matrix $G = (g_{ij})_{MN \times MN}$.

For images of fixed size M by N, every MN^{th} order an dpositive definite matrix G induces a euclidean distance. Calculating & comparing the Euclidean distance of database images from the test image recognizes the hand gesture.

D. Audio synthesizer

MATLAB supports an audio-player using an in-built function audioplayer by creating an audio object. It supports varies input arguments sampling rate, number of bits per sample for floating point signal whose valid values for are 8, 16, and 24 (by default it is 16) and identifier to specify the selected audio output device which is -1 for the default audio output device.

In our project we use pre-recorded audio for a particular hand gesture and by using the index of the database image the action is recognized

E. Implementation

The figure 7 shows the overview of the project,

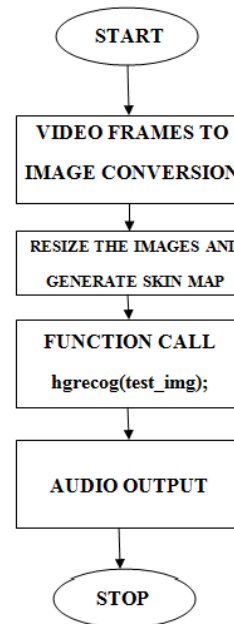


Figure 7 Overview of the project system

1) Video frames to image conversion

The step 1 of the implementation is video frames are converted to images by looping the frame2im() function, with one input argument the frame and imwrite() function with three input argument, the image to be written, the file name and format in which the image has to be written, for number of frame times individually by using two different for loops.

2) Resize the image

The step 2 of the implementation is to reduce the size of the images obtained from the step1 to a size of 280X280 in order reduce the memory required to store and process the images. MATLAB function imresize() which takes the input arguments as the image and the size in terms of number of rows and columns.

3) Generation of the skin map

The step 3 of the implementation is to read the image using imread() function and applying the Grey-world algorithm and applying RGB to YCbCr conversion, detecting the pixels that are within the skin tone range and replacing the intensity with white pixels intensity 255.

4) Hand gesture recognition function

In Step 4 of the algorithm the hand gesture recognition function is called with the input argument specifying the path of the query image.

The steps in this function are as follows:

Reorganise all two dimensional images in the database as training images into one dimensional column matrices. Then put these one dimensional column matrices in a row to build two dimensional matrix. Compute eigenHG, m and A to extract principle component analysis feature

- (1) m - (MxN)x1 average of the images in training database.
- (2) A - (MxN)xP Matrix of image vectors after each vector getting subtracted from the mean vector m.
- (3) eigenHG-(MxN)xP', P' Eigenvectors of Covariance matrix (C) of training database X, where P' is the number of eigen values of C that best represent the feature set.

For the considered [MxN] matrix, the highest count of eigen values with non-zero value that is the minimum of [M-1,N-1] is the size that its covariance matrix can posses.

Since the count of pixels of individual image vector is at peak in comparision with count of query images, hence count of eigen values with nonzero value of C will be max P-1 where P is the count of query images.

Compute eigen values and eigenvectors of $L = A'A$, with eigenvectors related to eigenvectors of C linearly.

Eigenvectors are computed from non-zero eigen values of C, the feature sets represented in such a manner is best. Kaiser's rule is used identify the eigenvector sthat is principle components to be considered.

If computed eigen value is more compared to 1, then the eigenvector will be taken for creation of the eigenHG.

$$\text{eigenHG} = A * L_eig_vec$$

PCA features are extracted for the query image by computing eigenHG, A, m of query image.

The comparison between the 2 gestures is done by projecting the gesture images onto gesture space and the Euclidean distance between them is measured.

Computing and comparison of the Euclidean distance of all projected test from the projected train images helps us to recognize the gesture.

5) Audio output

The step 5 of the algorithm uses the returned value of the hgrecof(test_img) function and by using the index as classifier the audio of the hand gesture is recognized. By using the function audioread() with input argument, the pre-

recorded audio to read an audio file along with the format followed by audioplayer() function with the input arguments as the audio file read and to create an audio object. The audio object created is passed as input argument to the play() function to obtain the audio output.

IV RESULTS AND DISCUSSIONS

The output to be considered here is voice output as the project aims towards it. In addition to that we also obtain intermediate results. Once the test images obtained by video are compared with the database images, we can see the recognized image of that particular action.

Light illumination, background colour, distance between hands and camera are some of the facts to be considered while video capturing. The proposed methodology was designed and tested with three set of actions. The actions considered are fire, yes, when respectively as shown in figure 8.

The input for feature extraction is the preprocessed gesture. Least Euclidean distance is computed between query and train images and gesture is recognized. Voice format output is obtained by the conversion of the recognized gesture.

The intermediate results we obtained will be resized image, image obtained after skin map generation and the images after comparison. We choose to show only the compared images in figure 9,10 and 11.

Once the correct compared images that is recognized and test image as shown in fig is obtained, we can get the corresponding voice output

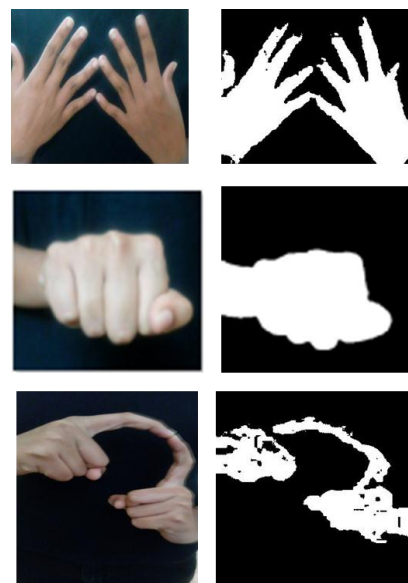


Figure 8 Actions for fire, yes, when respectively

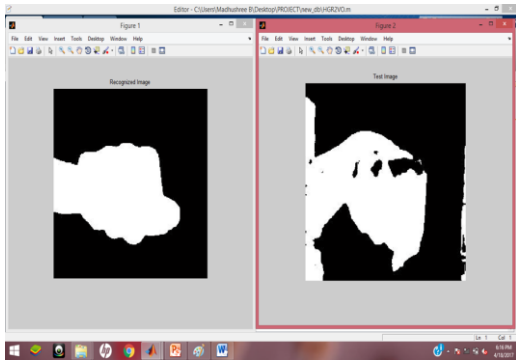


Figure 9. Output for recognition of yes

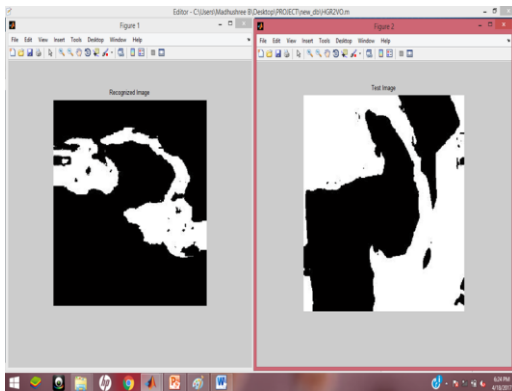


Figure 10. Output for recognition of when

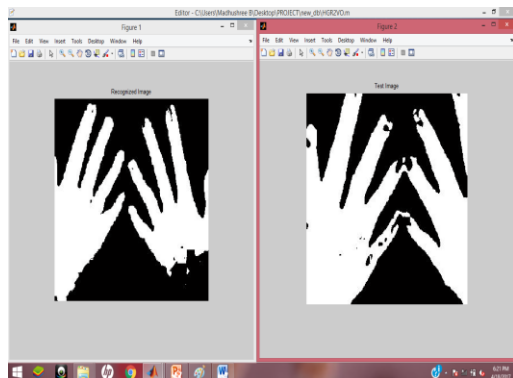


Figure 11. Output for recognition of fire

V CONCLUSION

An application executing the hand gesture recognition is Matlab based using Principle Component Analysis method was successfully implemented. The proposed technique gives text and audio output that aids to reduce the communication difference between mute & hearing impaired and normal people. Through this project, we have attempted to provide an artificial voice by recognizing the hand gesture. Action recognition can also be used for human to computer interaction. If we consider unfavorable and robust environments applicability will be more. We also need to deal with co-articulation that is accents of different people in different regions.

VI FUTURE SCOPE

The future scope of this work can be an apparatus that is developed as an aid for the people with no seeing ability. In this project the obstacles that are present in front of the user are captured using a camera. The user is given the information about the distance between the camera and the user or the presence of any defects in the path by the camera from the computing equipment and even the clear view of the obstacles if present in the path can be given with the help of audio synthesizer. We can also implement any sign language using this project. Further it can be improved to get output for combination of many words.

REFERENCE

- [1] Real-time handGesture detection and recognition using bag-of-features and support vector machine techniques. Nasser H. Dardas and Nicolas D. Georganas. IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT, 2011.
- [2] Visual recognition of american sign language using hidden markov models. Thad Eugene Starner. Master's thesis, Massachusetts Institute of Technology, Cambridge MA, 2015.
- [3] Hand-Gesture Recognition for Automated Speech Generation. Sunny Patel, Ujjayan Dhar, SurajGangwani, Rohit Lad, Pallavi Ahire. IEEE International Conference On Recent Trends In Electronics Information Communication Technology, May 20-21, 2016, India.
- [4] Static Vision Based Hand Gesture Recognition Using Principal Component Analysis. Mandeep Kaur Ahuja & Amardeep Singh. 3rd International IEEE Conference on MOOCs, Innovation and Technology in Education (MITE) 2015.
- [5] Hand Gesture Recognition of English Alphabets using Artificial Neural Network. Sourav Bhowmick, Sushant Kumar and Anurag Kumar. IEEE 2nd International Conference on Recent Trends in Information Systems (ReTIS) 2015.
- [6] Smart Glove With Gesture Recognition Ability For The Hearing And Speech Impaired. Tushar Chouhan, Ankit Panse, Anvesh Kumar Voon and S. M. Sameer. IEEE Global Humanitarian Technology Conference - South Asia Satellite (GHTC-SAS) September 26-27, 2014.
- [7] Sign Language Recognition. Anup Kumar, Karun Thankachan and Mevin M. Dominic. 3rd InCI Conf. on Recent Advances in Information Technology I RAIT-2016.
- [8] Real Time Sign Language Recognition using PCA. Shreyashi Narayan Sawant, M. S. Kumbhar. IEEE International Conference on Advanced Communication Control and Computing Technologies (ICACCCT) 2014.
- [9] Vision Based Hand Gesture Recognition Using Dynamic Time Warping for Indian Sign Language. Washef Ahmed, Kunal Chanda, Soma Mitra. International Conference on Information Science (ICIS) 2016.
- [10] Multiple Sign Language Translation into Voice Message. Hussana Johar R.B, Priyanka A, Revathi Amrut M S, Suchitha K, Sumana K J. International Journal of Engineering and Innovative Technology (IJEIT) Volume 3, Issue 10, April 2014.
- [11] Review in Sign Language Recognition Systems Symposium on Computer & Informatics. M. Ebrahim Al-Ahdal & Nooritawati Md Tahir, (ISCI), pp:52-57, IEEE ,2012.
- [12] Sign Language to Speech Converter Using Neural Networks. Mansi Gupta, Meha Garg, Prateek Dhawan. International Journal of Computer Science & Emerging Technologies 14 Volume 1, Issue 3, October 2010.
- [13] Embedded Based Hand Talk Assisting System for Deaf and Dumb. J. Thilagavathy, Dr.Sivanthi murugan, S. Darwin. International Journal of Engineering Research & Technology (IJERT), vol 3, issue 3, March 2014.
- [14] Hand Gesture Recognition Systems: A Survey. Arpita Ray Sarkar, G. Sanyal, and S. Majumder, International Journal of Computer Applications, vol. 71, no.15, pp. 0975 -8887, May 2013.