

A Vision Based Approach for Pakistan Sign Language Alphabets Recognition

NabeelSabir Khan (Corresponding Author)

School of Science and Technology (SST),

University of Management and Technology, Lahore

PO box 54770, Lahore, Pakistan

Tel: +9242111300200 E-mail: nabeel.bloch@umt.edu.pk

Adnan Shahzada

Politecnico di Milano, Milan, Italy

E-mail: shahzada@elet.polimi.it

Saleem Ata

School of Engineering (SEN),

University of Management and Technology, Lahore

PO box 54770, Lahore, Pakistan

Tel: +9242111300200 E-mail: saleemata@umt.edu.pk

Adnan Abid

School of Science and Technology (SST),

University of Management and Technology, Lahore

PO box 54770, Lahore, Pakistan

Tel: +9242111300200 E-mail: adnan.abid@umt.edu.pk

YaserDaanial Khan

School of Science and Technology (SST),

University of Management and Technology, Lahore

PO box 54770, Lahore, Pakistan

Tel: +9242111300200 E-mail: yaser.khan@umt.edu.pk

Muhammad ShoaibFarooq

School of Science and Technology (SST),

University of Management and Technology, Lahore

PO box 54770, Lahore, Pakistan

Tel: +9242111300200 E-mail: shoaib.farooq@umt.edu.pk

Dr. M. Tahir Mushtaq

School of Engineering (SEN),

University of Management and Technology, Lahore

PO box 54770, Lahore, Pakistan

Tel: +9242111300200 E-mail: tahir.mushtaq@umt.edu.pk

Inayat Khan

Centre of Excellence in Science and Applied Technologies, Islamabad, Pakistan

E-mail: inayatk@gmail.com

Abstract

There is a persistent communication barrier between the deaf and normal community because a normal person has no or limited fluency with the sign language. A person with hear-impairment has to express himself via interpreters or text writing. This inability to communicate effectively between the two groups affects their interpersonal relationships. There are about 0.24 million Pakistanis who are either deaf or mute and they communicate through Pakistan Sign Language (PSL). In this research work a system for recognizing hand gestures for Pakistan Sign Language alphabets in unimpeded environment is proposed. A digital camera is used to acquire PSL alphabet's images with random background. These images are preprocessed for hand detection using skin classification filter. The system uses discrete wavelet transform (DWT) for feature extraction. Artificial neural network (ANN) with backpropagation learning algorithm is employed to recognize the sign feature vectors. The dataset contains 500 samples of Pakistan Sign Language alphabets with various background environments. The experiments show that the classification accuracy of the proposed system for the selected PSL alphabets is 86.40%.

Keywords: Pakistan Sign Language (PSL), Discrete Wavelet Transform (DWT), Computer Vision, Artificial Neural Network (ANN), Skin Classification Filter, Back Propagation Learning Algorithm

1. Introduction

Dumb and deaf people face great difficulty while interacting with normal people. As the hearing impaired or deaf people cannot talk like normal people, so they are dependent on a kind of visual communication most of the times. A sign language is a way of communication that uses visual modality instead of acoustic sound patterns for exchanging information [1]. People with hearing and speech impairment use different sign and gestural languages in their daily face-to-face communication. There are different mutually unintelligible sign languages around the world, because the languages were developed independently of other deaf communities. American Sign Language, British Sign Language, French Sign Language, Australian Sign Language and many others have developed in various deaf communities. Similarly, Pakistan Sign Language (PSL) is developed in Pakistan by the deaf and dumb community which has its own vocabulary and syntax.

A sign language recognition system can be useful in bridging the communication gap between hearing impaired and the normal people. Because normal people are usually unaware of sign language grammars, it is very difficult for them to understand what a dumb or deaf person is trying to communicate. As a result, communication of a dumb person is usually limited only within the family or the deaf community [2]. A sign language recognition system may facilitate the deaf community by translating the sign language to natural

language. We propose a mechanism through which signs can be captured, recognized and translated to the text. Researchers have generally used either the direct devices or the computer vision based interfaces for Sign Language recognition systems. The direct devices include data gloves, styli and other position tracking devices[3], while vision based approaches use a camera to capture the hand gestures and movements. The main advantage of using computer vision based approach is that the user need not wear or use any complicated and expensive device.

Researchers have applied numerous techniques such as hidden markov model, principal component analysis, statistical measures and artificial neural networks. M. Atiqur Rahman et al. have used height, area, centroid, and distance of the centroid from the origin (top-left corner) of the image as features and then extracted features are used to train a Back propagation NN[4]. A. Giegal has used principal component analysis to recognize American Sign Language alphabets[5]. A neural network based speech translator is developed by Mansi Gupta et al. that undergoes the process of conversion from RGB to LAB to binary form. Then the number of black pixels for each block is computed and saved. After this the Euclidean distance between the signed input and stored data is compared for classification[6]. SignTutor is an interactive system proposed by Oya Aran et al. for Sign Language

Tutoring that uses a glove-based interface[7]. Ali Karamiet et al. have implemented a system that uses statistical measures for dimension reduction and ANN for classification of Persian sign language[8]. The system proposed by Ayoub Al-Hamadi et al. uses translation, rotation and scale invariant features for recognizing postures[9]. Gesture and sign language recognition has been the focus of many researchers but unfortunately there is no significant work done for Pakistan sign language at the moment. There is a system named "Boltay Hath" that aims at recognizing Pakistan Sign Language but it uses data gloves as its interface.[10]. The statistics given by Population Census Organization (Government of Pakistan), there are more than 3.3 million Pakistanis who are disabled. About 0.25 million among them suffer from hearing loss which is around 7.4% of the overall disabled population in the country[11]. A significant part of the deaf population is young and sign language recognition system can turn them into useful human resources for certain positions.

Our contribution in the paper is three fold: 1) The proposed system is the only computer vision based system for PSL recognition to date. 2) The system does not require any hardware like cyber-gloves that in effect reduces the cost and cumbersomeness of wearing hardware devices.

3) The proposed system works efficiently with random background gestures.

This research was conducted to propose an intelligent, robust and efficient system for the dumb and deaf people that will help them to communicate with other people in their gestural language. The proposed system recognizes the Pakistan sign language alphabets using image processing, discrete wavelet transforms and artificial neural networks. The organization of the paper is as follows: Section 2 explains Pakistan Sign Language's origin and alphabets. There is a brief introduction of Discrete Wavelet Transform in section 3 and section 4 explains the proposed system. Experimental results are discussed in section 5. The conclusion of whole research work is given in section 6.

2. Pakistan Sign Language

Pakistan Sign Language has its own vocabulary, syntax and semantics. It has undergone continuous improvement and evolution like all other languages. Spoken languages of an area have a significant impact on the growth and development of the sign language and variety of blends occur whenever they interact with each other. Signed English is a dialect of sign language that is formed by the combination

of British Sign Language and English. Sign Exact English (SEE) is gestural language that matches each spoken word of American English language. Likewise Signed Urdu has emerged by the combination of PSL, Urdu and other regional languages (Punjabi, Sindhi, Pushtu, Baluchi)[10].

Urdu Language consists of 38 alphabets but usually we find 37 unique hand gestures for most significant alphabets as shown in Fig. 1. Short messages are sent and received using these hand shapes. Hearing-impaired people in Pakistan make use of these hand shapes from PSL vocabulary combined with small gestures to express the words in Urdu. As Urdu is a combination of languages so English Sign alphabets

are also used with PSL. Forexample, to represent Saturday, 'S' is represented using both hands from English sign alphabets followed by the sign of Saturday. Such variations are bound to exist due to the existence of different cultural backgrounds in the same region through our history. Urdu itself is spoken in many different ways in different regions of Pakistan and it differs in vocabulary, phonology and grammar too. Similarly PSL has regional variations in many items. There are many examples that one sign is acceptable in one region but not preferred in another region in the same country but that does not conclude that the sign is wrong for that particular region.



Fig.1. Pakistan Sign Language Alphabets

3. Discrete Wavelet Transform (DWT)

Introduction of the wavelet transform was motivated by the need of further developments from Fourier Transform (FT). Although, FT gives us the information about the frequencies present in a signal but does not explain about the locality of the frequency components. The wavelet transform was introduced at the beginning of the 1980s by Morlet and has many applications in the areas as mathematics, physics, signal processing, medical imaging and image processing. The information provided by the continuous wavelet transform (CWT) is too redundant to reconstruct the signal and requires a lot of computations.

This is where discrete wavelet transform (DWT) comes into play that not only provides the sufficient information to analyze and reconstruct a signal but also reduces the computational requirements.

The foundations of the DWT were laid by Croiser, Esteban, and Galand by devising a technique to decompose discrete time signals back in 1976. Crochiere, Weber, and Flanagan also worked on the similar lines in that year.

They named their analysis scheme as subband coding. In the discrete case, filters of different cutoff frequencies are used to analyze the signal at different scales [15]. In DWT the signal is passed through iterations of low and high pass filters and then rescaling is done. Filtering is used to determine the signal resolution that represents the amount of detail information in the signal. Scaling on the other hand is achieved by up sampling and down sampling of the original signal [16], [17].

As image is a two dimensional (2D) signal, we can implement 2D discrete wavelet transform by applying low and high pass filters to decompose the image in different coefficients and down samplers to reduce the data without losing any of the information. The decomposition results in cutting down the spatial resolution in half since only half the number of samples now characterizes the entire signal but because frequency band of signal now contains only half of the previous bands, the frequency resolution is doubled.

This process can be repeated for further decomposition of the image. Each decomposition level will have half number of samples and double frequency resolution than the previous level. This 2D-DWT leads to a decomposition of approximation coefficients at some level in four components:

- The approximation Image: (Both horizontal and vertical direction trends)
- Vertical Detail Image: (Low frequencies in horizontal and high frequencies in vertical direction)
- Horizontal Detail Image: (High frequencies in horizontal and low frequencies in vertical direction)
- Diagonal detail image: (Detail image in both, horizontal and vertical directions).

The DWT of image $I(x,y)$ of size $M \times N$ is:

$$W_{\phi}(j_0, m, n) = \frac{1}{\sqrt{MN}} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} I(x, y) \phi_{j_0, m, n}(x, y) \quad (1)$$

$$W_{\phi}^i(j, m, n) = \frac{1}{\sqrt{MN}} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} I(x, y) \phi_{j, m, n}^i(x, y) \quad (2)$$

where j_0 is an arbitrary starting scale and $W_{\phi}^i(j, m, n)$ coefficients define an approximation of $I(x,y)$ at scale j_0 .

The $W_{\phi}^i(j, m, n)$ coefficients add horizontal, vertical and diagonal details for scales $j \geq j_0$ and $i = H, D, V$.

Low frequency components of the image constitute the approximation image whereas high frequency components make the three (Vertical, Horizontal and Diagonal) detailed images.

4. The Proposed System

The proposed system to recognize Pakistan Sign Language (PSL) alphabets can be divided into two major phases:

- Feature Extraction
- Classification

4.1. Feature Extraction

The process of Feature extraction consists of the following components:

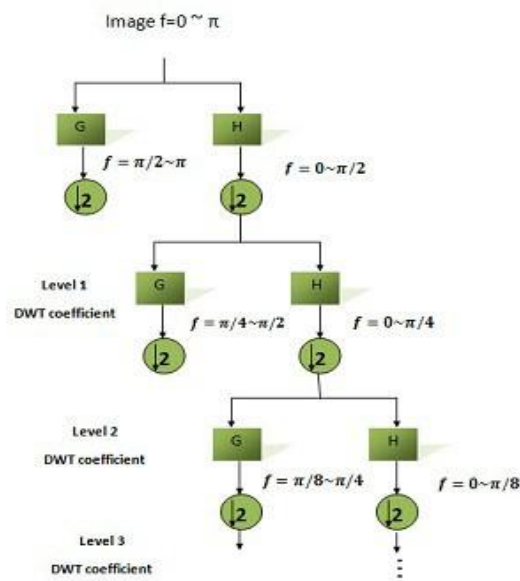


Fig.2. Discrete Wavelet Transform (DWT) based decomposition

4.1.1) Illumination and Color Balancing:

The system inputs RGB (Red, Green, Blue) images taken under normal lighting conditions using ordinary digital camera so they may have many inconsistencies and variations among them. In real world, the assumption that given an image with sufficient amount of color variations, the average value of the R, G, and B components of the image should average to a common gray value is held very well. Therefore, we can use this assumption to force the images in our dataset to have the same gray level value for each of the three (RGB) channels to reduce the effects of luminance inconsistencies. We have used Gray World algorithm [12] for color balancing and illumination compensation. It calculates the average values of R, G and B color components of the image and then determines scaling factors for all three components based on the deviations they have with the common gray value of the image.

4.1.2) Hand Segmentation:

The next step is to extract the hand out of the image for which color based segmentation is used. The resultant RGB image from the Gray World Algorithm is then converted to YCbCr (Y is luminance and both Cb and Cr are the Chromatic components) color space [13]. YCbCr form is derived by from the corresponding RGB space as follows:

$$Y = 0.2989R + 0.5866G + 0.1145B \quad (3)$$

$$Cr = 0.7132(R - Y) \quad (4)$$

$$Cb = 0.5647(B - Y) \quad (5)$$

The following Chromatic components of the YCbCr color space image are used to detect pixels that appear to be skin using the range of Cb and Cr that Chai and Ngan [14] have found to be the representative of the skin color.

$$77 \leq Cb \leq 127$$

$$133 \leq Cr \leq 173$$

The skin pixels are marked as blue by setting the pixel[R G B] values to [0 0 255] correspondingly. It is then converted to a binary image by setting all the skin pixels as 1 and non-skin pixel as 0. A 4x4 median filter is applied then to fix the isolated false classified skin pixels. The Figure.3 shows the result of segmentation.



Fig.3.ColorBaseSegmentation

4.1.3) Cropping and Resizing:

Cropping is done by eliminating all the zero valued rows and columns from the binary images and then the images are resized to 300x400 pixels as shown in the Figure. 4.



Fig.4.Cropping and Resizing

4.1.4) Discrete Wavelet transform:

Our system makes use of DWT property of retaining the same information component despite of less number of samples to reduce the dimensions and extract features out of an image. Haar wavelet transform was used to derive the interest points from the sign images because it shows better results while describing the human body parts [18]. As we are using artificial neural networks for the classification, it is not possible to use all the image elements (pixels) as the input. Therefore, we have used 2D-DWT at multiple levels to reduce the number of inputs to the neural network. We have experimented with different levels of decomposition through DWT to find the best coefficient matrices that may produce the coefficients nearest to the original one when reconstructed. In this paper, we have used the coefficients of approximation along with the detail coefficients on the 6th level. Feature extraction through DWT can be visualized by the Figure. 5



Fig5. Discrete Wavelet Transform (DWT): Approximation, Horizontal, Vertical and Diagonal detail images

4.2. Classification

Artificial neural network (ANN) is used as the signlanguage alphabets classification model in our proposed system. ANN is an information processing system that takes its inspiration from biological neural network. Artificial neural network aims to generalize the mathematical models of neural biology and imitates the human learning and cognitive processes. Number of constituent neuron layers and the pattern of connections between them is offered as the Architecture of the network. The way connection strengths (weights) are adjusted is called learning algorithm and activation function determines and controls the output of the neurons against the set of inputs [19].

4.2.1) ANN Architecture:

A feed forward Multilayer Perceptron (MLP) is used for the classification of alphabets. We used a single hidden layer because we did not find any substantial benefit of adding another hidden layer. In case of neural network, you have to empirically find the right architecture for the classification task at hand. We have run many simulations with different network structures and finally found the network with 140 input neurons, a single hidden layer with 75 neurons and 37 output neurons exhibiting the best results. See Figure. 6.

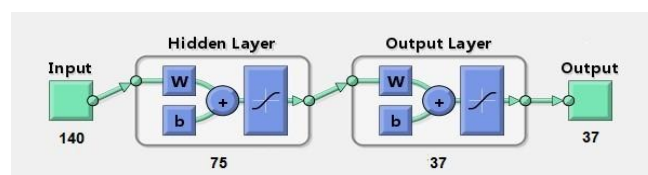


Fig.6. Artificial Neural Network (ANN) Structure

4.2.2) Training Algorithm and Parameterization:

Gradient descent with momentum and adaptive learning rate backpropagation is used to train the neural network. This algorithm works very well with the noisy and inconsistent data and improves generalization [20]. Tan-Sigmoid is used as the activation function for all the layers.

$$T(S) = \tanh(S) = \frac{e^S - e^{-S}}{e^S + e^{-S}} \quad (6)$$

It generates the output in the range [-1,1].

Training is done in supervised manner and a training sample is a pair $((x_p, d_p))$ where $p = 1, \dots, P$. The input vectors are denoted by $(x_p = x_{p,1}, \dots, x_{p,n})$, output is represented by $(o_p = o_{p,1}, \dots, o_{p,n})$ and desired output is $(d_p = d_{p,1}, \dots, d_{p,n})$. Error in the classification for a single pattern is defined as:

$$E_{p,j} = o_{p,j} - d_{p,j} \quad (7)$$

The back-propagation algorithm applies a correction $\Delta W_{a,b}$ to the synaptic weight $W_{a,b}$ which is proportional to the partial derivative $\frac{\partial E_{p,j}}{\partial W_{a,b}}$ that is called gradient descent.

$$\frac{\partial E_{p,j}}{\partial W_{a,b}} = -E_j T_j(S) o_{p,j} \quad (8)$$

The performance of the network is measured by Mean Squared Error (MSE) defined as the average squared error between the network generated outputs and the target outputs.

$$MSE = \sum_{p=1}^P \sum_{j=1}^K E_{p,j}^2 \quad (9)$$

Random data division is used for validation.

5. Experiments and Simulation Results

A digital camera was used to capture 37 static onehandedPSL alphabets. The experimental data was collected by varying the hand orientation, its distance and angle with the camera with random backgrounds. Our dataset contained 500 images of 37 alphabets from which 426 images were utilized for training and 74 for testing. The proposed neural network described in the previous section was simulated using Matlab. The performance curve and train state of the ANN is shown in Fig. 7 and 8.

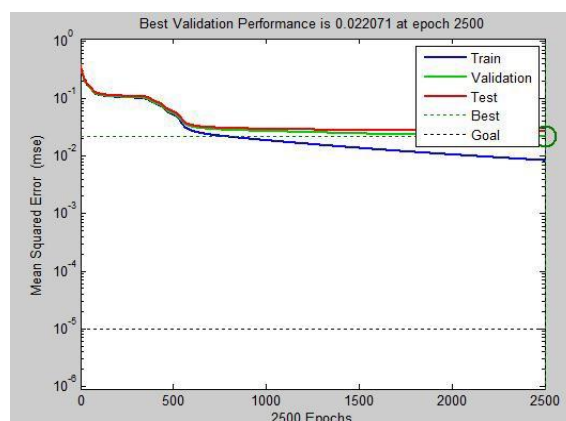


Fig.7.Performance Curve during Training Process

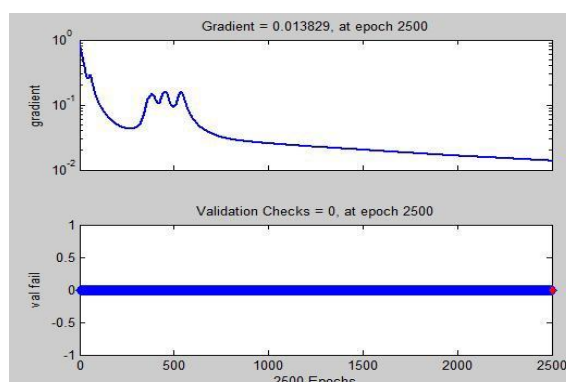


Fig.8.Training State

The experimental results for both the training and testing data are shown in Table. 1.

Table 1: Experimental Results Representation

Data	Total Samples	Correctly Classified	Accuracy (%)
Training	426	365	85.10
Testing	74	63	84.08
Overall	500	425	84.6

It can be observed that our recognition system was able to classify data with 84.60% even with the random background. Further, Figure. 9 shows the accuracy with which all individual alphabets are recognized by the system.

Finally the system is compared to the models suggested by "BoltayHaath". These models include Statistical Interval Matching (SIM), Statistical Interval Matching Combined with LMS and Statistical Interval

MatchingCombined with Democracy "[10]. Our proposed system beats or matches the accuracy of all of the three models mentioned. The result comparisons are shown in Table. 2

Table 2: Comparison with BoltayHaath

Model	Accuracy (%)
Statistical Interval Matching	26
SIM Combined with LMS	84
SIM Combined with Democracy	73
Our Proposed System	84.6

Our results demonstrate that the proposed model hashigh generalization capability and performs efficiently forvaried background environments.

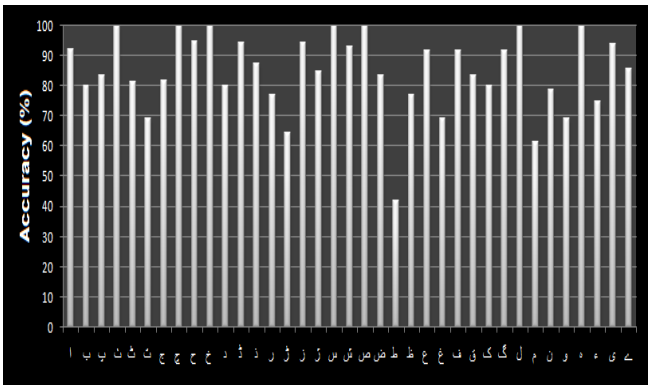


Fig.9.Representation of Classification Results of PSL Alphabets

6. Conclusion

An intelligent system for gesture and sign recognition isproposed in this paper to facilitate a better communicationamong normal and the deaf communities. The proposedsystem recognizes the Pakistan Sign Language alphabetswith random backgrounds. A digital camera is used toget images of PSL alphabets instead of data gloves.

The system segments the hand from the images usingskin color tracking and extracts the features by applyingdiscrete wavelet transform (DWT) in the first phase. Inthe second phase, the extracted features are applied toNeural network for classification. We used dataset of 500samples of Pakistan sign Language alphabets in variousbackground environments and the experiments show thatthe proposed system is able to classify the selected PSLsigns with a classification accuracy of 84.60%. The networkis trained using MATLAB NN Toolbox.

References

- [1] Armstrong, D.F., Stoke, W.C., Wilcox, S.E., "Gesture and nature of language", Cambridge Academic Press, 1995.
- [2] Foez M. Rahim, Tamnun EMursalin, Nasrin Sultana. "Intelligent Sign Language Verification System, Using Image Processing, Clustering and Neural Network Concepts". International Journal of Mathematical Engineering (IJME), Vol. 1, No. 1-2, Jan-Dec, 2009, pp 43-56
- [3] M. Mohandes, S. A-Buraiky, T. Halawani and S. Al-Baiyat. "Automation of the Arabic Sign Language Recognition". Information and Communication Technologies: From Theory to Applications, 2004. Proceedings. 2004 International Conference on. On page(s): 479-480
- [4] Md. Atiqur Rahman, Ahsan-Ul-Ambia and Md. Aktaruzzaman "Recognition Static Hand Gestures of Alphabet in ASL". 2011 IJCIT, VOLUME 02, ISSUE 01.
- [5] A. Geigel, "Recognizing American Sign Language Letters Using Principal Component Analysis". CPSC 320 Vision IRC/November 2004
- [6] Mansi Gupta, Meha Garg, Prateek Dhawan, "Sign Language to Speech Converter Using Neural Networks", International Journal of Computer Science and Emerging Technologies (E-ISSN: 2044-6004) Volume 1, Issue 3, October 2010
- [7] Oya Aran, Ismail Ari, "Alexandre Benoit Sign Tutor: An Interactive System for Sign Language Tutoring". IEEE MultiMedia, January-March 2009, pp. 81-93
- [8] Ali Karami, Bahman Zanj, Azadeh Kiani Sarkaleh, "Persian sign language (PSL) recognition using wavelet transform and neural networks". Expert Systems with Applications: An International Journal Volume 38 Issue 3, March, 2011
- [9] Ayoub Al-Hamadi, Omer Rashid and Bernd Michaelis, "Posture Recognition using Combined Statistical and Geometrical Feature Vectors based on SVM". International Journal of Information and Mathematical Sciences 6: 12010.
- [10] Aleem Khalid Alvi, M. Yousuf Bin Azhar, Mehmood Usman, Suleman Mumtaz, Sameer Rafiq, Razi Ur Rehman, Israr Ahmed. "Pakistan Sign Language Recognition Using Statistical Template Matching". World Academy of Science, Engineering and Technology volume 3, 2005
- [11] Population Census Organization, Govt. Of Pakistan at <http://www.census.gov.pk/index.php>
- [12] G.D. Finlayson, B. Shiele and J.L. Crowley. "Comprehensive colour normalization". Proc. of European Conference on Computer Vision, (ECCV). Vol. I, 475-490, Freiburg, Germany, 1998.

- [13] Charles Poynton, *"Digital Video and HDTV"*, Chapter 24, pp. 291-292, Morgan Kaufmann 2003.
- [14] D. Chai, and K. N. Ngan, *"Face segmentation using skin-color map in videophone applications"*. IEEE Trans. on Circuits and Systems for Video Technology, 9(4):551-564, June 1999.
- [15] Robi Polikar, Tutorial on: *"Multiresolution Analysis: The Discrete Wavelet Transform"*
- [16] Mallat, S. (1989), *"A theory for multiresolution signal decomposition: the wavelet representation"*; IEEE Pattern Anal. and Machine Intell., Vol. 11, pp. 674-693
- [17] Misiti, M., Misiti, Y., Oppenheim, G., & Poggi, J.-M. (2006). *"Wavelet toolbox for use with MATLAB"*. <<http://www.mathworks.com>>.
- [18] Dinh, T. B., Dang, V. B., Duong, D. A., Nguyen, T. T., and Le, D. *"Hand gesture classification using boosted cascade of classifiers"*. In 2006 International conference on research, innovation and vision for the future. pp. 139-144.
- [19] Haykin, S. (1999). *"Neural networks: A comprehensive foundation"*, (2nd ed.). New Jersey: Prentice Hall.
- [20] Munib, Q., Habeeb, M., Tahruri, B., and Al-Malik, H. A. (2007). American sign language (ASL) recognition based on Hough transform and neural networks. Expert Systems with Applications, pp. 24-37.