# Gesture Recognition Based on Hexagonal Structure Histograms of Oriented Gradients

Pang Haibo, Liu Chengming, Zhao Zhe and Zhang Shuyan

*School of Software Technology, Zhengzhou University, Zhengzhou, China, 450002*
*phb@zzu.edu.cn*

## *Abstract*

*Feature extraction methods of image directly affect the feature recognition results based on computer vision-based gesture recognition. This paper proposed a gesture feature extraction method of hexagonal structure, which is based on histograms of oriented gradients. This paper transforms the quadrilateral structure into hexagonal structure of the image and defines hexagonal block structure. Using different structures, different sizes of hexagonal blocks extract gesture feature and recognize on different sizes image. Experimental results show that hexagonal structure block is more appropriate than conventional histograms of oriented gradients structure block, the hexagonal structure histograms of oriented gradients feature extraction method is more efficient than the square and circular structure histograms of oriented gradients feature, using gentle_adaboost classifier achieved higher recognition rate.*

*Keywords: Feature extraction; Hexagonal structure; Histograms of oriented gradients; Gesture recognition*

## 1. Introduction

In recent years, there has been great interest in the studies related to vision-based HCIS (Human Machine Interaction Systems). Such as face recognition, body-pose recognition, hand-gesture recognition, voice recognition, fingerprint recognition, and so on [1].

Hand-gesture recognition has become an important application in HCIS, because keyboard and mouse cannot satisfy people's interacting requirements to be great extent. Up to now, there have been many successful algorithms for hand-gesture recognition, but there are still some factors affecting the performance of the hand-gesture recognition, say, illumination, masking, etc. Therefore, the main tasks for hand-gesture recognition are to find better methods for features extraction and robust classification.

In general, gesture recognition is considered as a very challenging old since natural environments tend to be rather unsuitable for gesture recognition due to bad illumination, nonuniform backgrounds, and so on. The numerous publications of the recent years show that static hand gesture recognition is still an old of active research, whereas many of them try to face the previously mentioned problems in order to improve the performance and quality of existing technologies.

The remaining part of this paper is organized as follows: Section 2 introduces the related works. Section 3 describes our approach, including the details of hexagonal structure histograms of oriented gradients. In section 4, we introduce how to acquire hand-gesture images, extract h-hog features. Furthermore, we present the experimental processes and results. Conclusions and suggestions for our future works are given in Section 5.

## 2. Summary of the Related Works

Chen and Tseng [2] proposed a multi-angle hand-gesture recognition method in their article. They used three webcams set at, front, right, left directions of hand to capture gestures. Then three SVM (Support Vector Machine) classifiers trained respectively. After the training process, one voting and two plans of fusion fused the constructed classifiers. The recognition rate of hand gestures more than *93*%, including different angles, sizes, and different skin colors. However, there were only three hand gestures in their research.

Huang and Hu [3] applied the method of gabor filter, PCA and SVM to recognize many simple hand gestures in complex background. They first extracted the hands from a sequence of video images using the skin color information. Then, they coped with these images of hand gestures using gabor filter, and they used PCA method to reduce the dimension of the data space and used SVM classifier to recognize the hand gestures. The recognition rate of *95.2*% can be achieved. The experiment results show that the processing time is *0.2* second for every frame, which did not achieved the requirement of a real-time system.

Huang and Hu [4] estimate the orientation of the hand gestures using the gabor filter responses. The estimated angle is used to correct the hand pose into an upright orientation.

Amin and Yan [5] proposed a system that is able to recognize ASL (American Sign Language alphabets) from hand gesture with average *93.23*%. They used PCA and gabor filters to accomplish the task, out of the top *20* principal components the best combination of principal components is determined by finding the best fuzzy cluster for the corresponding PCs of the training data, their experiment demonstrate the best result obtained from the combination of the fourth to seventh principle components. However, their recognition rate of similar alphabets is relatively low.

Jayashree R [6] uses fixed position low-cost web camera with *10* mega pixel resolution mounted on the top of monitor of computer which captures snapshot using RGB color space from fixed distance. This gesture recognition system can reliably recognize single-hand gestures in real time and can achieve a *90.19*% recognition rate in complex background with a "minimum-possible constraints" approach.

Li [7] proposed a speed hand gesture recognition system using the kinect sensor. Based on the HOG features and adaboost training algorithm, the experiment demonstrated the detection results is great, but for the situations like hands covered in front of body or objects kind of similar to hands, there is still some high missing and false rate. The application of hand gesture recognition for real-life is very challenging because of the requirements on its robustness, accuracy and efficiency.

Padam Priyal [8] has presented a gesture recognition system using geometry based normalizations and Krawtchouk moment features for classifying static hand gestures. The proposed system is robust to similarity transformations and projective variations.

Cao [9] and Zhang [10] using image for depth extraction and recognition of shape matching method in hand gesture recognition has obtained the certain effect, but susceptible to the influence of occlusion.

The literatures [11-12] combined with histogram of oriented gradient and local gradient direction for recognizing gesture, and achieved good results. The literature [13] identify the hotspot function of the workload on an embedded system that motivates acceleration and present the detailed design of a hardware accelerator for histograms of oriented gradients descriptor extraction.

The literature [14] incorporated fuzzy concept to HOG aiming to achieve a good recognition rate with a low feature vector dimension. Experimental results have demonstrated that HFOG outperforms the original HOG with a lower dimensional vector. But HOG structure affected the classification performance.

# 3. Hexagonal Structure Histograms of Oriented Gradients

## 3.1 Histograms of Oriented Gradients

HOG was put forward earliest by Dalal [15], and initially to be used for human detection and achieved good results. It is based on the direction of the shape, appearance and other features of an image can be described by gradient or edge of the density distribution in principle, and it also can adapt to illumination changes and rotation of objectives.

Unlike other image geometric features, HOG does not consider the image characteristics from a whole, but subdivided image into several small cell units, and then calculate the gradient or edge direction histogram of each pixel in all cells. To improve performance, a number of cells form an block, thus the image become a connected graph consisting of some blocks, then normalize the gradient of all cells in these blocks, obtained the final direction of gradient vector.

The gradient direction vector of the image is decided by the size of image, the size of cell, and the dimension of gradient from each cell. The calculation method is as below:

$$Num = cDim * \left(\frac{bSize}{cSize}\right) * \frac{(w - bSize + bStep) * (h - bSize + bStep)}{bStep^2} \tag{1}$$

Where *cDim*, *bSize*, *cSize*, *bStep* denote cell dimensions, block size, cell size and block step respectively; *w* and *h* represent image width and height, *Num* is the number of HOG features. As *Num* is relatively large, in practice, typically need to adjust these parameters.

According to the actual situation in the practical application, we adjust the picture size, cell dimensions and other parameters, and strive to reach a suitable vector dimension, so that the amount of computation will not be too large for machine learning and hand recognition.

However, the edge direction information, and brightness of the original HOG structure block cannot be extracted to achieve satisfactory results when were applied to human detection and classification. Therefore, Dalal [16] and Zhu [17] improved some deformation Rectangular-HOG and Circular-HOG. Figure 1 shown some instances.
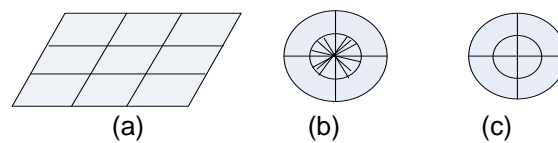


(a)          (b)          (c)

**Figure 1. Variants of HOG Descriptors**

## 3.2 Virtual Hexagonal Structure

We transfer the original image into hexagonal grid image, and then extract the hexagonal block structure feature.
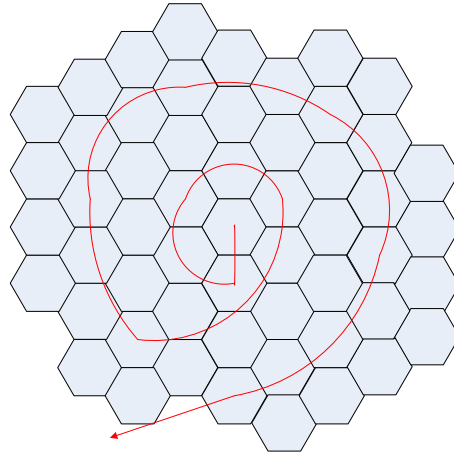
**Figure 2. Spiral Architecture with Spiral Addressing**

The possibility of using a hexagonal grid to represent digital images has been studied for more than thirty years. Hexagonal grids have higher degrees of symmetry than the square grids. This symmetry results a considerable saving of both storage and computation time. Figure 2 shown a one-dimensional addressing scheme for a hexagonal structure, called spiral architecture. From figure 2, it is easy to see that the location (denoted by $L$) of the pixel with a given spiral address：

$$a_n a_{n-1}...a_1, \quad a_i = 0,1,2...,6, \quad i = 1,2,...,n \tag{2}$$

$a_i$ can be found from the locations of (3):

$$a_i \times 10^{i-1}, i = 1,2,...,n \tag{3}$$

Using (4) ：

$$L(a_n a_{n-1}...a_1) = \sum_{i=1}^{n} L(a_i \times 10^{i-1}) \tag{4}$$

Because there has been no hardware available for image display and capture on hexagonal structure. To construct hexagonal pixels, each square pixel is first separated into $7 \times 7$ small pixels, called sub-pixels. To be simple, the light intensity for each of these sub-pixels is set to be the same as that of the pixel from which the sub-pixels are separated. Each virtual hexagonal pixel is formed by *56* sub-pixels as shown in figure 3. To be simple, the light intensity of each constructed hexagonal pixel is computed as the average of the intensities of the *56* sub-pixels forming the hexagonal pixel. Fig. 3 shows a collection of seven hexagonal pixels constructed with spiral addresses from *0* to *6*.

Let us assume that original images are represented on a square structure arranged as *2M* rows and *2N* columns, where *M* and *N* are two positive integers. Let the centre of the virtual hexagonal structure be located at the middle of rows *M* and *M+1*, and at column *N*. Note that there are 1*4M* rows and *14N* columns in the (virtual square) structure consisting of virtual sub-pixels obtained from the original square pixels. Let us construct the first hexagonal pixel using the *56* sub-pixels with centre located in the middle of rows *7M* and *7M+1* and the column *7N* of the virtual square structure.

After the *56* sub-pixels for the first hexagonal pixel are allocated, all sub-pixels for all hexagonal pixels can be assigned from formula (4).
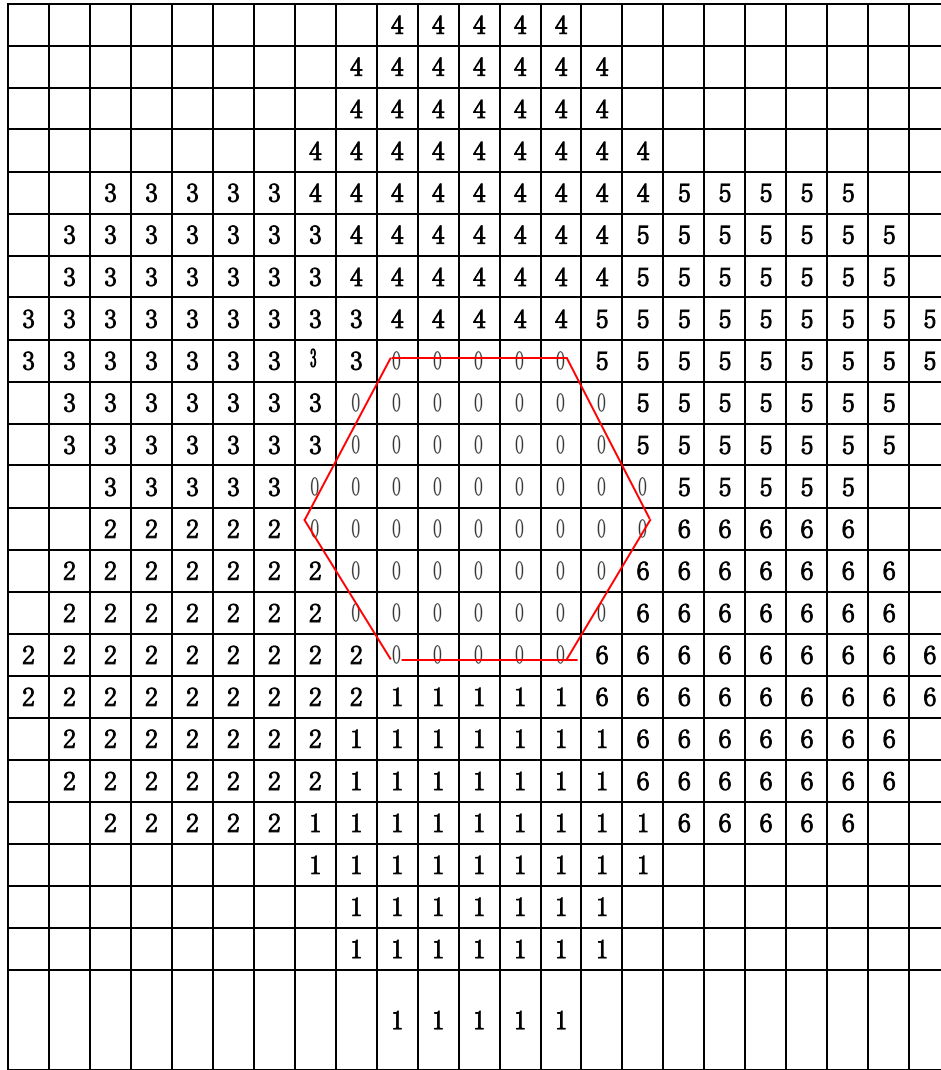
| | | | | | | | | | 4 | 4 | 4 | 4 | 4 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | 4 | 4 | 4 | 4 | 4 | 4 | 4 | | | | | | | |
| | | | | | | | | | 4 | 4 | 4 | 4 | 4 | 4 | 4 | | | | | | | |
| | | | | | | | | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | | | | | | |
| | 3 | 3 | 3 | 3 | 3 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 5 | 5 | 5 | 5 | 5 | | | |
| | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | |
| | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | |
| 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 4 | 4 | 4 | 4 | 4 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| | | 3 | 3 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 5 | 5 | 5 | 5 | | |
| | | 2 | 2 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 6 | 6 | 6 | | |
| | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 |
| | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | |
| | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | |
| | | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 6 | 6 | 6 | 6 | 6 | | |
| | | | | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | | | | | | | | | |
| | | | | 1 | 1 | 1 | 1 | 1 | 1 | 1 | | | | | | | | | | | | |
| | | | | 1 | 1 | 1 | 1 | 1 | 1 | 1 | | | | | | | | | | | | |
| | | | | | 1 | 1 | 1 | 1 | 1 | | | | | | | | | | | | | |

**Figure 3. A Cluster of Seven Hexagonal Pixels**

## 3.3 Extracting H-HOG Features

Basing on above-mentioned transformation from quadrilateral structure to hexagonal structure, the descriptor of hexagonal structure histogram of oriented gradient was proposed and shown in figure 4.

Similar to the C-HOG, H-HOG is a series of concentric hexagons whose center pixel are $d_i$, and radiuses of the circumscribed circles are $R$. The H-HOG definition is shown in formula (5).

$$h_{d_i}^0(k) = \#\left\{ d_j \neq d_i : d_j \in E_o : (d_j - d_i) \in bin(k) \right\} \quad k=1,...,(n_r \cdot n_\theta) \qquad (5)$$

The shape context of point $d_j$ is a two-dimensional distribution. $K_r$ and $K_\theta$ is respectively radial serial number and angle serial number of the $k$th container. The value of $K_r$ and $K_\theta$ is the numbers of pixel within the range of the triangle. "#" represents statistics number of points at triangle edges from which the distance is $K_r$ away to $d_i$, angle is $K_\theta$ and gradient is $o$.

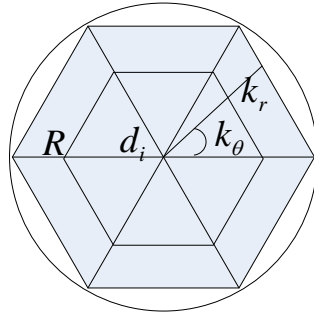$$k_r = R \sin 60^0 / \sin(120^0 - k_\theta) \qquad (6)$$

**Figure 4. H-HOG Descriptor**

If $r_j$ and $\theta_j$ represent the radial distance and angle of line from point $d_j$ to point $d_i$, generalization container fields of $d_j$ be defined as formula (7):

$$S(d_j) = \left\{ bin(l), l \in \left\{ (\theta_{d_j} + i) \bmod n_\theta \right\}_{i=-b}^{b} \right\} \tag{7}$$

Equation (5) can be rewritten as:

$$h_{d_i}^0(k) = \#\left\{ d_j \neq d_i : d_j \in E_o : (d_j - d_i) \in S(d_j) \right\} \tag{8}$$

$$SC_{d_i}^o(k) = \sum_{d_j} g(\theta_{d_j}, k) \cdot h_{d_i}^o(k), k = 1, ..., (n_r \cdot n_\theta) \tag{9}$$

$g(\theta_{d_j}, k)$ represents a Gaussian density function whose center is $\theta_{d_j}$ and peak value is *1*。

$SC_{d_i}^o(k)$ is the shape context of $u_j$.



**Figure 5. H-HOG Block on the Hexagonal Structure Image**

Figure 5 shows the result when a structure block diagram of the Figure 4 is applied to the Figure 2.

Figure 5 shows the result after the structure block of Figure 4 applied to the diagram of Figure 2. The solid lines show the three different scales of the hexagonal block, which include *7* pixels, *19* pixels and *37* pixels respectively. The dotted line shows a hexagonal block where contain *19* pixels also. The solid line block and the dotted line block have the overlapping pixels which contains the same *10* pixels in Figure 5.

# 4. Experiments and Results

## 4.1 Acquiring and Preprocessing Gesture Images

The training and testing hand-gesture images for our experiment were collected by capturing videos of different hand gestures. We recorded *12* hand-gesture videos with *320*240* pixels, which were extracted to frames, then we chose *50* images respectively *12* hand gestures (see Figure 6), including different angles (see Figure 7). Therefore, we obtained *50* images of each of hand gesture and *600* images for training and testing.

In order to speed up the experiments and verify the classification capability, we converted the hand-gesture images into grayscale images and resized to *40*30* pixels.



**Figure 6. Hand-gesture Images of Figure**



**Figure 7. Different Angle of Figure "1" and "3"**



**Figure 8. Hand-gesture Image of Figure 3**

Figure 8 is the original gesture image. The upper left corner gesture image of fig 9 is the approximate gesture image of third level wavelet composition which the resolution is *40\*30* pixels.
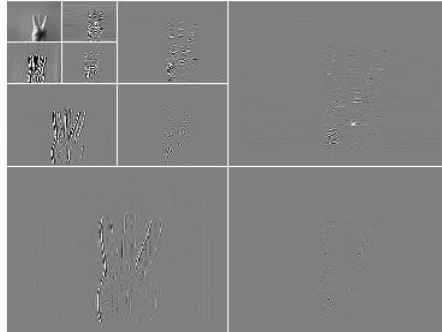


**Figure 9. The Third Level Wavelet Decomposition of Figure 8**

## 4.2. Designing Experimental and Analyzing Results

Experimental environment: Inter (R) core (TM) 2 Quad CPU Q8300@2.5GHz 2.5GHz, 2.00 GB memory.

Gentle_adaboost classifier was classified hand-gesture images. The experimental given the ROC (Receiver Operating Characteristic Curve) graphs and the recognition rates.

In our experiments, the dataset consist of *600* hand-gesture images, we randomly divided their features into training set (including *450* images) and testing set (including *150* images) in each experiment. The static gesture images have the resolution *320\*240* and *40\*30*, then extract R-HOG, C-HOG and H-HOG features. Finally, we introduced the gentle_adaboost classifier into our experiments, which "cross-validation" is the earliest and one of the most widely used implementations, and performed the classification operation.

The procedure was repeated *20* times, and then the average AUC (Area Under roc Curve) value and the recognition rate were counted.

The hexagon block was divided into *6* units which are assigned *0.2*, *0.2*, *0.2*, *0.2*, *0.1* and *0.1* respectively anticlockwise from the top right corner according to Figure 4.

Parameters include winSize (*64,128*), blockSize (*16*, *16*), blockStride (*8*, *8*), cellSize (*8*, *8*), nbins (*9*) when extract R-HOG feature. Meanwhile, parameters include four angular bins, *2* radial bins, *4* pixels radius for the center bin, *2* pixels expansion factor for the radius when extract C-HOG feature.

The HOG descriptor was the vector of the components of the normalized cell histograms from all of the block regions. These blocks typically overlap, meaning that each cell contributed more than once to the final descriptor.

Let *v* be the non-normalized vector containing all histograms in a given block; then the normalization factor can be the following formula:

$$f = v \Big/ \sqrt{\| v \|_2^2 + e^2} \tag{10}$$

To improve the recognition efficiency, we introduced the PCA to reduce the features dimension to *50*-dimensional.

Fig 10 showed the average TPR (True Positive Rate) and the average FPR (False Positive Rate) of testing set corresponding to R-HOG feature, C-HOG feature and H-HOG feature where hexagon block edge length was *3\*3* pixels and step was *2* pixel which images size were *40\*30*. According to the figure, we have known the AUCs were *0.9456*, *0.9623* and

*0.9537* respectively. Meantime, we could see that the recognition rates were *0.9221*, *0.9294* and *0.9283* from Figure 12.

Figure 11 showed the average TPR and FPR of testing set corresponding to R-HOG feature, C-HOG feature and H-HOG feature where hexagon block edge length was *3\*3* pixels and step was *2* pixel which original images size were *320\*240*. According to the figure, we have known the AUCs were *0.9538*, *0.9643* and *0.9771*. Meantime, the recognition rates were *0.9271*, *0.9386* and *0.9459*. Comparing to the two results, we had a conclusion that scale will affect the classification performance and the recognition rates.



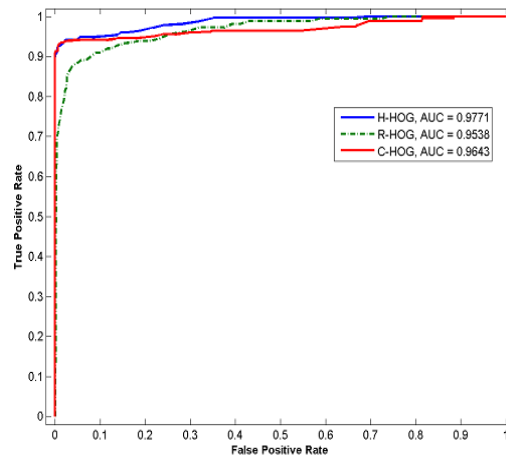**Figure 10. The R-HOG, C-HOG and H-HOG ROC when Image Scale is 40\*30 Pixels**



**Figure 11. The R-HOG, C-HOG and H-HOG ROC when Image Scale is 320\*240 Pixels**

In our experiments, we had used *12* different kinds of gesture; meanwhile, some gesture images were very similarly. We got a recognition rate of *97.83%*.

Figure 13 shows two different angle gesture-images of number "10". Due to rotation of hand gesture, it was difficult to divide them into the same type of gesture.

Figure 14 show similar gesture images of number "1" and "9". Due to rotation of hand gesture, the two gestures were so similar that cannot be distinguished.
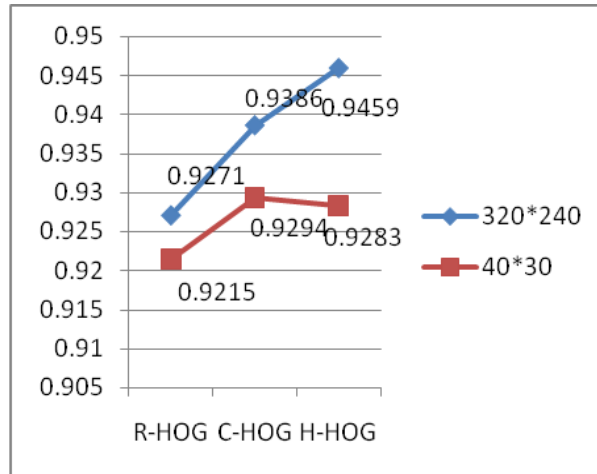
**Figure 12. The R-HOG, C-HOG and H-HOG Recognition Rates when Image Scale is 40*30 and 320*240 Pixels**



**Figure 13. Two Angle Gesture Images of Figure 10**



**Figure 14. Similar Gesture Images of Figure 1 and 9**

This article also test the performances of H-HOG structural block at conditions of side block length *3 * 3*, step *1*, block side length *4 * 4*, step *2*, block side length *4 * 4*, step *3*, block side length *5 * 5*, step *2*, block side length *5 * 5*, step *3*, block side length *5 * 5*, step *4*. The table 1 showed the efficiency and superiority of our descriptor.

The results of H-HOG are better than C-HOG and R-HOG, and we can draw the two following conclusions. First, our method improves over the C-HOG and R-HOG, especially for some large resolution gesture images. Second, our method does not significantly improve accuracy on other larger dataset.

**Table 1. H-HOG Accuracies when Image is 40*30 and 320*240 Pixels**

|         | 3*3(1) | 3*3(2) | 4*4(2) | 4*4(3) | 5*5(2) | 5*5(3) | 5*5(4) |
|---------|--------|--------|--------|--------|--------|--------|--------|
| **40*30**   | 0.9103 | 0.9294 | 0.9301 | 0.9421 | 0.9206 | 0.9301 | 0.9267 |
| **320*240** | 0.9312 | 0.9459 | 0.9630 | 0.9783 | 0.9551 | 0.9680 | 0.9593 |

## 5. Conclusion

In our article, we had proposed an algorithm of hand-gesture recognition based on H-HOG descriptor and Gentle_adaboost classifier, effectively addressing the static hand-gesture recognition problem. At first, the gesture images were transforming into hexagon images. Then we obtained H-HOG features when hexagon block edge length and step length were taken different values. Further, we trained and tested H-HOG features. At last, we compared R-HOG, C-HOG and H-HOG features performance.

Additionally, the recognition time fit to the requirement of a real-time system. Further, it is robust to nonlinear illumination and image blurring.

At present, 12 hand gestures had been used in our experiment. Some hand-gesture images were very similar; therefore, the classification error rate of our method was relatively high. Further research needed to research new feature descriptor.

## Acknowledgments

## References

[1]   S. P. Priyal and P. K. Bora, "A study on static hand gesture recognition using moments [C]", In: Proceedings of the International Conference on Signal Processing and Communications (SPCOM), (**2010**), pp. 1–5.

[2]   Y.-T. Chen and K.-T. Tseng, "Multiple-angle hand gesture recognition by fusing SVM classifiers [C]", In IEEE conference on Automation Science and Engineering, Scottsdale, AZ, USA, (**2007**) September, pp. 527-530.

[3]   D.-Y. Huang, W.-C. Hu, and S.-H. Chang, "Vision-based hand gesture recognition using PCA + Gabor filters and SVM[C]", In proceedings of the 5th International Conference on Intelligent Information Hiding and Multimedia Signal Processing, Kyoto, Japan, (**2009**) September, pp. 1-4.

[4]   D.-Y. Huang, W.-C. Hub and S.-H. Chang, "Gabor filter-based hand pose angle estimation for hand gesture recognition under varying illumination [J]", Expert Systems with Applications, vol. 38, (**2011**), pp. 6031–6042.

[5]   M. A. Amin and H. Yan, "Sign language finger alphabet recognition from Gabor-PCA representation of hand gestures [C]", In Proceedings of the Sixth International Conference on Machine Learning and Cybernetics, Hong Kong, (**2007**), August, pp. 2218-2223.

[6]   J. R. Pansare, S. H. Gawande and M. Ingle, "Real-Time Static Hand Gesture Recognition for American Sign Language (ASL) in Complex Background [J]", Journal of Signal and Information Processing, vol. 3, (**2012**), pp. 364-367.

[7]   H. Li, L. Yang, X. Wu, S. Xu and Y. Wang, "Static Hand Gesture Recognition Based on HOG with Kinect [C]", 2012 4th International Conference on Intelligent Human-Machine Systems and Cybernetics. 2012, pp. 271-273.

[8]   S. PadamPriyal and P. K. Bora, "A robust static hand gesture recognition system using geometry based normalizations and Krawtchouk moments [J]", Pattern Recognition, vol. 46, (**2013**), pp. 2202–2219.

[9]   C. Cao, R. Li and L. Zhao, "Hand Posture Recognition Method Based on Depth Image Technology [J]", Computer Engineering, vol. 38, no. 8, (**2012**), p. 4.

[10]  P. Zhang, T. Li, H. Xiong and L. Liang, "Gesture Recognition Based on Depth Difference Distribution [C]", 21st International Conference on Pattern Recognition (ICPR 2012), pp. 157-160, November 11-15, 2012. Tsukuba, Japan.

[11]  M. B. Kaaˆniche and F. Breˊmond, "Gesture Recognition by Learning Local Motion Signatures [C]", IEEE Conf. Computer Vision and Pattern Recognition, (**2010**).

[12]  M.-B. Kaaˆ niche and F. Breˊmond, "Recognizing Gestures by Learning Local Motion Signatures of HOG Descriptors [J]", IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 34, no. 11, (**2012**) November, pp. 2247-2258.

[13] S. E. Leea, K. Minb and T. Suhc, "Accelerating Histograms of Oriented Gradients descriptor extraction for pedestrian recognition [J]", Computers & Electrical Engineering, vol. 39, no. 4, (**2013**) May, pp. 1043–1048.

[14] A. I. Salhi, M. Kardouchi and N. Belacel, "Histograms of fuzzy oriented gradients for face recognition [C]", 2013 International Conference on Computer Applications Technology, (**2013**) January 20-22, pp. 1-5.

[15] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection [C]", In IEEE Conference on Computer Vision and Pattern Recognition, (**2005**) June, pp. 886-893.

[16] N. Dalal, "Finding People in Images and Videos [D]", The French National Institute for Research in Computer Science and Control, Grenoble, France, (**2006**).

[17] Q. Zhu, S. Avidan and M. C. Yeh, "Fast Human Detection Using a Cascade of Histograms of Oriented Gradients [C]", In: Proceeding of IEEE Conference on Computer Vision and Pattern Recognition, (**2006**), pp. 1491-1498.

# Authors

**Pang Haibo**, he was born in 1979, PH.D, and Lecturer. His research interests include image processing and pattern recognition.

**Liu Chengming**, he was born in 1979, PH.D, Lecturer. His research interests include image processing.

**Zhao Zhe**, he was born in 1983, Master, and Lecturer. Her research interests include image processing and virtual reality.

**Zhang Shuyan**, he was born in 1987, Master, and Lecturer. Her research interests include image processing and intelligent algorithm.