

## 지휘행동 이해를 위한 손동작 인식

제홍모<sup>o</sup>, 김지만, 김대진

포항공대

[invu71@postech.ac.kr](mailto:invu71@postech.ac.kr), [jmk@postech.ac.kr](mailto:jmk@postech.ac.kr), [dkim@postech.ac.kr](mailto:dkim@postech.ac.kr)

### Hand Gesture Recognition for Understanding Conducting Action

Hongmo Je<sup>o</sup>, Jiman Kim, Daijin Kim

POSTECH

#### ABSTRACT

We introduce a vision-based hand gesture recognition for understanding musical time and patterns without extra special devices. We suggest a simple and reliable vision-based hand gesture recognition having two features : First, the motion-direction code is proposed, which is a quantized code for motion directions. Second, the conducting feature point (CFP) where the point of sudden motion changes is also proposed. The proposed hand gesture recognition system extracts the human hand region by segmenting the depth information generated by stereo matching of image sequences. And then, it follows the motion of the center of the gravity(COG) of the extracted hand region and generates the gesture features such as CFP and the direction-code. Finally, we obtain the current timing pattern of beat and tempo of the playing music. The experimental results on the test data set show that the musical time pattern and tempo recognition rate is over 86.42% for the motion histogram matching, and 79.75% for the CFP tracking only.

#### 1. Introduction

Many reports on intelligent human machine interaction using hand gesture recognition have already been presented [1]. Without specialized tracking devices, one of the greatest challenges of the system is to reliably detect and track the position of the hands using computer vision techniques. The vision-based hand gesture recognition methods use only the vision sensor; camera [2]. In general, the entire system of the vision-based hand gesture recognition must be more simple than the Data Glove-based approach, and it makes human-friendly interaction with no extra device. The vision-based hand gesture recognition is a challenging problem in the field of computer vision and pattern analysis, since it has some difficulties of algorithmic problems such as camera calibration, image segmentation, feature extraction, and so on.

A vision-based method for understanding human's conducting action for chorus with a special purpose light baton and infrared camera have been suggested in [3]. It proposed a vision system which captures the image sequences, tracks each end-point of the baton which is a stick having a distinguished color

feature to be detected easily by a camera, and analyzes a conducting action by fuzzy-logic based inference. Lately, Watanabe and Yachida [4] have proposed a real-time interactive virtual conducting system using the Principle Component Analysis (PCA)-based gesture recognition that can identify only 3/4 time pattern.

In general, conductors perform various music using both hands and natural conducting actions may be very difficult to represent. Hence, we take a few assumptions to make the problem easy as follows :

- 1) The conductor uses only one-side hand
- 2) The conducting action must be in the view range of the camera.
- 3) the conductor may indicate four timing patterns (2/4, 3/4, 4/4, 6/8) with three tempos (Andante, Moderato, Allegro) by his/her hand motion.
- 4) the conductor needs no special devices.

We propose a very simple but reliable vision-based hand gesture recognition of the human conductor with no extra devices. Unlike the previous vision-based hand gesture

recognition, we use the depth information, instead of using the intensity or the color information of image, generated by a stereo vision camera to extract human hand region that is the key region of interest (ROI) in this application. Our proposed system can obtain the motion velocity and the direction by tracking the center of gravity (COG) of the hand region, which provides the speed of any conducting time pattern. We introduce two methods to recognize the musical time pattern. One is the *CFP tracking* which uses only special features like conducting feature point and another is the *motion histogram matching* which can identify the time pattern and the tempo at once, where the "Mahalanobis distance" is chosen as the distance metric of motion histogram matching.

The remainder of this paper is organized as follows. Section 2 describes the proposed hand gesture recognition system to understand music time pattern and tempo in detail. Section 3 presents the experimental results of both simulation and real world videos. Finally, Section 4 draws conclusion and discusses future work.

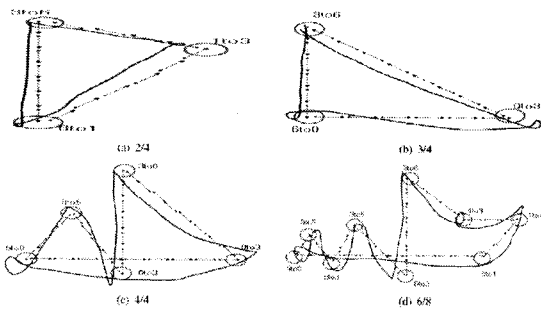


Fig.1. The real trajectories and the approximated directions of each conducting pattern. (solid line - real trajectory, dashed line - motion direction, red circle - conducting feature point).

## 2. The Proposed Hand Gesture Recognition System

The system has two stages which are 'hand segmentation' and 'music time pattern and tempo recognition'. More details on each stage are described in following subsections.

### 2.1 Hand Segmentation Stage

Hand segmentation separates the human hand region from the others. Most methods for the hand segmentation uses the skin color information to extract the hand region. The skin color-based hand region extraction is quiet simple, but it is sensitive to light condition change and complicated and cluttered background which has many skin-like colored objects such as wood and wall papers. We use the depth information of a stereo image instead of the 2D pixel image. The depth information might not only be insensitive in any light condition but also

robust even if there is a complicated background. We utilize a face detector to detect the human face, which allow us to find the hand candidate region easily because we know that the depth of hand region must be closer than the face region to the stereo camera.

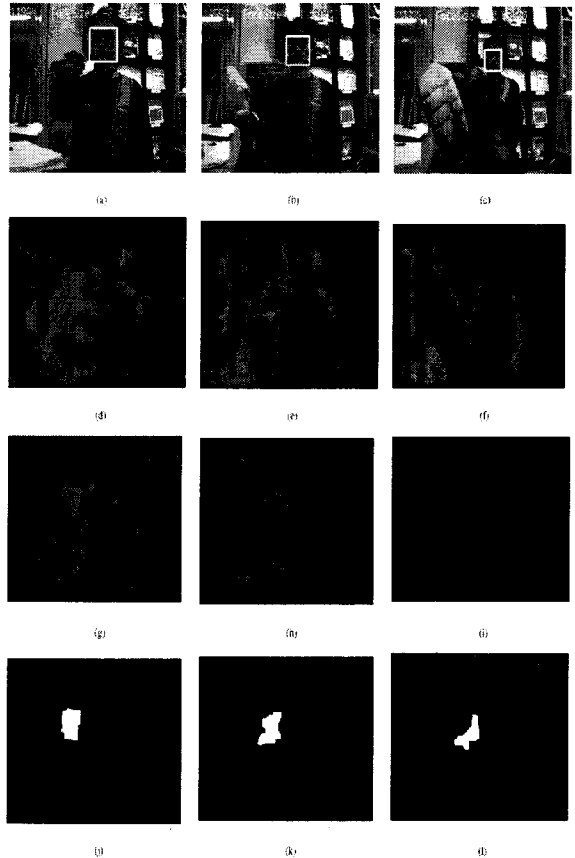


Fig. 2 (a),(d), (g), (j) A human conductor stands at 1.0m from the camera, (b),(e), (h), (k) 1.5m from the camera, (c),(f), (i), (l) 2.0m from the camera. ; the first row (a), (b), (c) show images ; the second row (d), (e), (f) show the depth maps ; the third row (g), (h), (i) show the noise removal and segmentation; the fourth row (j), (k) (l) show the result after connected component analysis and morphological operation.

Fig. 2 shows the intermediate results through hand segmentation stage. The exact hand region can be completely segmented after postprocessing on the extracted hand candidate region. Fig. 2 shows several postprocessing stages such as the morphological operator and the connected component analysis to remove the non-hand region [5]. To track the motion of the hand, the COG of the hand region needs to be computed. We approximate the COG by computing the mean coordinates of the segmented hand region as

$$X_{cog} = \frac{\sum x_i}{N}, Y_{cog} = \frac{\sum y_i}{N} \quad (1)$$

, where  $x_i, y_i$  are the x and y coordinates at the  $i$ th pixel position, respectively, and  $N$  is the number of pixels of the hand region.

**2.2 Musical Time Pattern and Tempo Recognition**

In contrast to the hand sign language recognition, the hand gesture recognition for understanding a musical time pattern and tempo does not have to be accurate. While a slight posture or movement of hands in the hand sign language represents an independent and important meaning, only salient features like beat transition point of hand motion are the most important information in the conducting gesture.

**The direction code of the hand motion**

The easiest way to find the trajectory of the conducting gesture is to track the motion direction of the hand. We obtain the direction angle of the hand motion by computing the difference between the previous COG of hand region and the current COG of it as

$$\begin{aligned} \Delta X_{cog}(t) &= X_{cog}(t) - X_{cog}(t-1), \\ \Delta Y_{cog}(t) &= Y_{cog}(t) - Y_{cog}(t-1), \\ \theta(t) &= \arctan \frac{\Delta Y_{cog}(t)}{\Delta X_{cog}(t)}, \end{aligned} \quad (2)$$

where  $\theta(t)$  is the direction of the hand movement on time  $t$ .

To represent the direction-code, the real value of the hand direction should be quantized in eight directions. Fig. 3 shows three-bit codes for the eight dominant direction of hand movement.

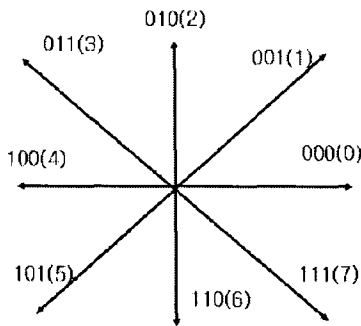


Fig. 3. Eight direction codes.

**Conducting feature point**

Fig. 1 illustrates the representative features which are called as

“Conducting Feature Point (CFP)”, of each musical time pattern. For example, a musical time pattern of 2/4 has three CFPs which are 3to6, 6to1, and 1to3. Assuming that the new coming CFP is 6to1 while the previous CFP is 3to6 or the start point of the initial gesture, then the gesture recognition system expects the next CFGs 1to3 following 3to6. Thus, the recognition system can identify the time pattern by observing the sequence of CFPs.

TABLE I. The profile index for the time patterns and tempos.

<b>Index</b>	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>
Time pattern	2/4	2/4	2/4	3/4	3/4	3/4
Tempo	And.	Mod.	Alle.	And.	Mod.	Alle.
<b>Index</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>	<b>11</b>
Time pattern	4/4	4/4	4/4	6/8	6/8	6/8
Tempo	And.	Mod.	Alle.	And.	Mod.	Alle.

**Motion histogram matching**

Although the analysis of CFP sequences is reliable for identifying musical time patterns, it can fail when the system misses an important CFP. This can occur for a complicated time pattern like 6/8 which has a large variation among the different human conductors. To avoid this problem, we propose a motion histogram matching based on the musical time pattern tempo analysis. We can obtain a cycle of each time pattern, where one cycle means a sequence of the direction-code from the start point of time patterns to their end point (usually both the start point and the end point are the same). In general, most musical time patterns have “3to6” type of the CFP as the start point of their action. In the training stage, we collect the histogram vectors  $H = [h_0, h_1, \dots, h_7]$  where  $h_0$  to  $h_7$  are the number of each direction code for the cycle and obtain the statistics (mean, variance) of motion histogram for all combinations of time patterns (2/4, 3/4, 4/4, 6/8) and tempos (Andante, Moderato, Allegro). Then, the mean and variance vectors  $H_\mu$  and  $H_\Sigma$  of the histogram vectors can be computed as

$$\begin{aligned} H_\mu &= \frac{1}{N} \sum_i H_i, \\ H_\Sigma &= \frac{1}{(N-1)} \sum_i (H_i - H_\mu)^2, \end{aligned} \quad (3)$$

where  $N$  is the number of training data. Thus, we have twelve profiles of motion histogram. Table I. denotes the profile

index for the time patterns and tempos. For example,  $H_{\mu}^1$  represents the mean of 2/4 with moder3to tempo and  $H_{\Sigma}^{11}$  represents the variance of 6/8 with allegro tempo. We selected the "Mahalanobis distance" as a metric of motion histogram matching. By Eq.(4), the similarity scores for all profile are evaluated. The proposed musical time pattern and tempo recognition system identify the time pattern and tempo by taking the profile whose similarity score is the minimum.

$$MD^k = \sqrt{(H_c - H_{\mu}^k)^T H_{\Sigma}^{k-1} (H_c - H_{\mu}^k)} \quad (4)$$

$$ProfileIndex = \arg \min_k MD^k, \quad (5)$$

where  $H_c$  is the current motion histogram and  $k$  is the profile index given in Table. I.

### 3. Experimental Results

We used the 'BumbleBee stereo vision camera' [6] as input sensor. Since it automatically provides depth information for each frame of stereo images, we needed not perform stereo matching with heavy computational complexity. We collected they conducting gesture data for each time pattern and tempo which consists of 300 cycles respectively. We divided them into 200 cycles for training the recognition system and 100 cycles for testing the recognition system. Fig. 4 represents the recognition rate of the experiments for the profile indices. As a result, the average recognition rate of using the CFP sequence analysis is 79.75% and that of using the motion histogram matching is 86.42%.

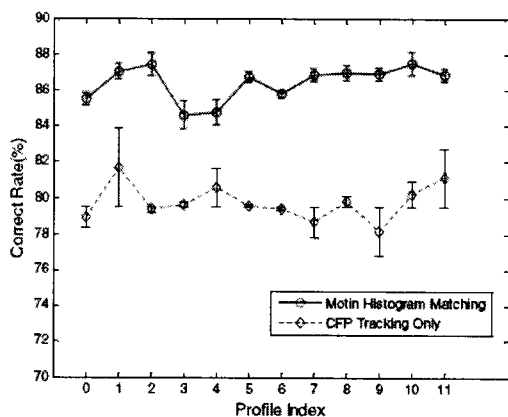


Fig. 4. The plot of recognition rate ; red solid line for 'Motion Histogram Matching', blue dashed line for 'CFP tracking only'

### 4. Conclusion

This paper presented a vision-based hand gesture recognition. We implemented a system for understanding musical time pattern and tempo that was generated by a human conductor. We only used the stereo vision camera with no extra devices. Instead of using the color pixel image, we used the depth information of the current stereo image. It is easy to extract an interesting region on the depth map. In addition face detection helped us to find the hand candidate regions because we assumed that the hand motion of the conductor usually is made in front of the human face. With this assumption, we could extract the hand region faster and easier. We introduced the conducting feature points with the direction-codes which simply indicated the musical time pattern and tempo that was generated by hand motions. When the human conductor made a conducting gesture, our proposed system tracked the COG of the hand region and encoded the motion information into the direction-code. We also suggested the motion histogram matching which could identify the current musical time pattern and tempo simultaneously by finding the best matched distribution of direction-code in a cycle. From the numerous experiments, the recognition accuracies were 79.75% and 86.42% using the CFP sequence analysis and the motion histogram matching, respectively.

#### ACKNOWLEDGMENT

This work is financially supported by the Ministry of Education and Human Resources Development(MOE), the Ministry of Commerce, Industry and Energy(MOCIE) and the Ministry of Labor(MOLAB) through the fostering project of the Lab of Excellency, and partially supported by the Korea Science and Engineering Foundation (KOSEF) through the Biometrics Engineering Research Center (BERC) at Yonsei University.

#### REFERENCES

- [1] Ying Wu and Thomas S Huang, "Vision-based gesture recognition: A review," *LNCS: Gesture-Based Communication in Human-Computer Interaction: International Gesture Workshop*, vol. 1739, pp. 103, 2004.
- [2] A. Mulder, "Hand gestures for hci," *Technical Report 96-1*, vol. Simon Fraster University, 1996.
- [3] Zeungnam Bien and Jong-Sung Kim, "On-line analysis of music conductor's two-dimensional motion," San Diego, CA, USA, 1992, pp. 1047-1053.
- [4] Takahiro Watanabe and Masahiko Yachida, "Real-time gesture recognition using eigenspace from multi-input image sequences," *IEEE Computer and System in Japan*, vol. 30, no. 13, pp. 810-821, 1999.
- [5] R. Bloem, H. N. Gabow, and F. Somenzi, "An algorithm for strongly connected component analysis in  $n \log n$  symbolic steps," pp. 37-54, Nov 2000, LNCS 1954.
- [6] "http://www.ptgrey.com/product"