

# Gesture Recognition Based on Kinect

Yunda Liu, Min Dong\*, Sheng Bi, Dakui Gao, Yuan Jing, Lan Li

School of Computer Science & Engineering  
South China University of Technology  
Guangzhou, China  
hollymin@scut.edu.cn

**Abstract**—With the rapid development of computer science, gesture recognition has been a highlight of research in the area of Human Computer Interaction (HCI). Generally speaking, gesture recognition can be divided into two types: static gesture recognition and dynamic gesture recognition. Under the background of the aging population, how to use the method of gesture recognition to help the elderly adapt to “Intelligent Age” is a meaningful issue, which deserves more attention. This paper describes a gesture recognition method based on Kinect, a 3D somatosensory camera sensor. This method involves skeleton tracking, where the skeleton data is produced from the depth images obtained via Kinect. Extensive experiments demonstrate the superior performance of the proposed methods over Kinect.

**Keywords**— *Human Computer Interaction; Kinect; gesture recognition; skeleton tracking; depth image*

## I. INTRODUCTION

Gesture is one of the natural ways of Human-Computer Interaction (HCI) and it has intuitive control methods. In general, gestures are defined as motions involving hands or the combination of hands and arms. It is usually divided into two types: static and dynamic gestures. Static gesture means the shape of single hand, which is corresponding to the relevant points in the model-parameter space. While dynamic gesture is composed of a series of motions, corresponding to a trajectory in parameter space, which is described by the space characteristics changing with time. Algorithms on static gesture recognition have developed rapidly in recent years, such as gesture recognition based on artificial neural network and computer vision. However, simple static gesture cannot meet the requirements of the industry application and methods on dynamic gesture recognition have become a focus in the area of research instead. Owing to the diversity and complexity of gestures and flexible positions and shapes of hands, it makes gesture recognition become a multi-discipline and challenging project<sup>[1][2][3][4]</sup>.

Traditional hand gesture recognition research began in the end of the 1980s. C. Charaphayan et al<sup>[5]</sup> proposed a static recognition using image-processing algorithm based on the words in ASL (American Sign Language). This method can accurately identify the 27 of the 31 words in the sign language. T. Starner et al<sup>[6]</sup> proposed the recognition of ASL by extracting the features of the images as a hidden Markov model. G. Bradski et al<sup>[7]</sup> came up a method about recognizing gestures for video content navigation. M. J. Jones et al<sup>[8]</sup> proposed a method based on the colors of skin. T. Starner et al<sup>[9]</sup> also proposed a gesture recognition system using two

different camera combinations. It is composed of skin color model and HMM and uses them to perform motion tracking and gesture training.

Since 2010 Microsoft issued Kinect, a kind of somatosensory camera, a number of methods based on it have been proposed. The former gesture system cannot identify and locate hand positions, but for Kinect, all these problems can be solved easily.

Kinect can obtain three kinds of data: RGB data, Depth data, and sound and presage data. Developers can get suitable data to develop their own application according to demands. There are two main methods based on Kinect for gesture recognition. The first one is to obtain human depth images, and separate the region of hands by the threshold value. The second one is to get human skeleton images using skeleton tracking through Kinect, then it uses the skeleton data for further gesture recognition.

In this paper, an effective gesture recognition system based on Kinect is proposed. This system combines the key points of skeletal trails matching with static key frames matching to carry on dynamic gesture recognition. And the remainder of this paper is organized as follows. In Section 2, principles of Kinect are introduced. Section 3 presents the recognition algorithm suggested in this paper. In Section 4, the experimental results about the algorithm are provided. The conclusion and future work are given in section 5.

## II. RELATED WORK

### A. Static Gesture Recognition

As a significant research field of computer vision, static gesture recognition usually applies algorithms such as template matching, machine learning and some approaches based on intelligent hardware. Generally speaking, the traditional procedure of static gesture recognition can be summarized into the following parts: separation of hand regions, hand tracking, description of hand characteristics and hand classification. With the advent of Kinect, it becomes more convenient to research and develop the algorithms of gesture recognition. Equipped with the color and depth cameras, Kinect can provide RGB as well as depth images, which are helpful to collect related data. To be mentioned, since depth images are not sensitive to light, humidity or other environmental factors of background, by combining them with RGB images, separation of the hand regions becomes more accurate and efficient. Once the job of separation is

completed more accurately, many other algorithms used afterwards for special matrix extraction or classification becomes more meaningful. For instance, in the next step, we could use the descriptor to extract features of the images, which are composed of a training data set. Then, many classic pattern recognition methods can be applied in the training process, such as Support Vector Machine (SVM), Decision Tree etc. Eventually, we classify the input static gestures from the test set and evaluate the accuracy of recognition for further improvement. When the accuracy of gesture recognition reaches a high level, this algorithm for static gesture recognition is deemed to be valuable.

### B. Principles of Kinect

#### (1) Generation of Skeleton Data

Skeleton tracking in Kinect is efficient, which can track and distinguish 20 main joints of human body. There are three steps for mapping from a depth image to its corresponding skeleton image.

Step1: Recognition on Human Body. Margin detection, noise threshold handling and classification of objective characteristics are used to separate the body from background. “Divide and conquer” is made for every tracked figure, which removes the environment image in the background.

Step2: Recognition on Body Parts. The result of first step is used and the human-depth image is generated without background. Every part of human body is recognized, extracted by fast recognition based on Kinect.

Step3: Recognition on Joints. Joint recognition involves machine learning to analyze the recognized human body, where the data of 20 joints is obtained. Considering that there will be overlapped or covered parts, Kinect analyzes and learns from three different views: front, side (left/right) and above. There are three states for every joint: tracked, not tracked and not inferred.

#### (2) Skeletal Space Coordinate

When users come into the vision of Kinect, the sensor can obtain 20 joints of human body. Therefore, when a user moves in the vision of Kinect, it starts to collect skeletal data and judge whether the motion makes any sense.

Kinect uses Cartesian Coordinate to define the position of joint in the vision, and the unit is meter. Actually users are in a skeletal space coordinate system whose origin is the Kinect itself. The coordinate system accords with right-hand rule, where x-axis is on horizontal direction. When users are facing Kinect, the direction of x-axis is defined as right-hand side, and the y-axis is vertical upward, the z-axis is toward the direction of Kinect facing the users.

#### (3) Skeletal tracking

Skeletal tracking contains two kinds of patterns: active and passive patterns. Kinect can track at most two persons, which means it can offer at most two persons' skeleton data. Our system mainly uses skeleton tracking under the active pattern.

According to datasets obtained by active pattern, the following information is achieved:

1) Active pattern contains the whole skeleton data corresponding to skeletal tracking, including spatial position information of the current user's 20 joints.

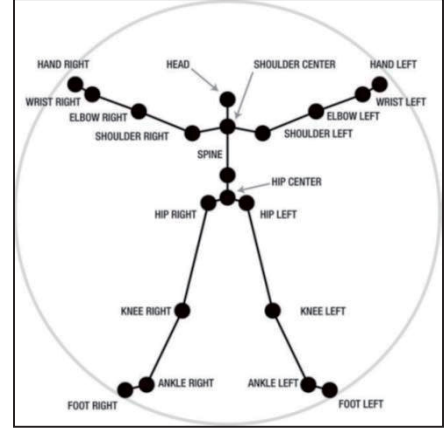


Fig. 1. Joint recognized by Kinect

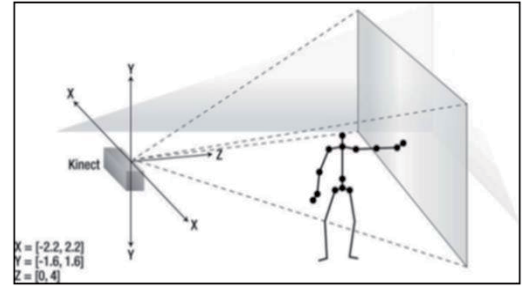


Fig. 2. Skeletal space coordinate

The unique skeletal tracking ID is distributed to every user in the field of view to distinguish them.

### III. DYNAMIC GESTURE RECOGNITION

In this session, the details of recognition methods are clarified. There are three basic dynamic gestures to be recognized: sliding left hand towards right, sliding right hand towards left and lifting up right hand. The main procedures of this method includes obtaining skeleton data, determining the left or right hand, specifying the relative parameters, error controlling and the display of the recognition results.

#### A. Obtaining skeleton data

Kinect randomly assigns IDs to users in the tracking zone and it is not effective to track a specific user. Besides, if users leave and then enter the tracking zone again, the returned index values may also change. In order to avoid these shortcomings, the person who is nearest to the sensor is defined as the target in the paper. After identifying the target, the skeleton data is added to a linked list. When the new skeleton data comes, it will be appended at the end of the list. If the list has been already full, it will be emptied. In order to get a more precise result, all joint (including hand joint) data should be collected.

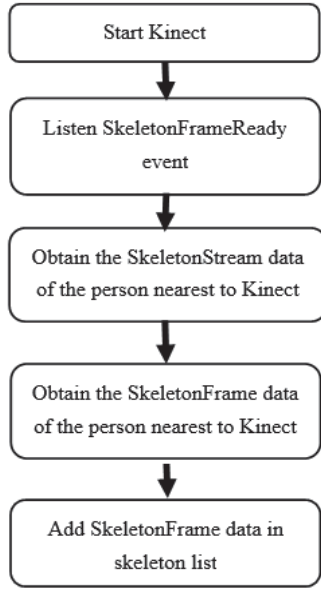


Fig. 3. Obtain skeleton data

#### B. Process of Gesture Recognition

The recognition of mentioned three gestures are based on the data stored in skeleton linked list as well as the following factors: error distance in vertical direction, direction of movement, duration of gestures completion, and the distance of movement.

When the gesture is performed, there will be slight changes in vertical direction. In our method, the differences of Y coordinates between two sequent joint data should be calculated to determine the gesture is done horizontally. Next, we should test in which direction the gesture is done. The differences of X coordinates between every two sequent data in the list. The positive result indicates the user moves his hand in one direction consistently. The gesture should also be completed within a range of time. A minimum distance value is involved to make sure the gesture is actually happening.

**Input:** The gesture performed by the user

**Output:** The result of the recognition process

Initializing error in vertical direction, time range, and minimum distance.

start=0

for index from 1 to the length of skeleton linked list

//Determine which hand is moving

if(jointID==HandLeft)

if list[index+1].X - list[index].X <0

start=index

if(jointID==HandRight)

if list[index+1].X - list[index].X >0

start=index

//Determine the hand is moving horizontally

if list[index].Y - list[start].Y >error in vertical direction

start=index

//Determine the gesture is happening

if list[index].X-list[start].X>minimum distance and duration is in the range

return true

return false

Algorithm. 1. Recognition for left hand waving to the right

In this algorithm, the size of skeleton list is 20. Because Kinect can capture data of 30 frames per second. The maximum threshold value of time to recognize is set as 0.75s. In this case, it needs 20 frames to complete the recognition.

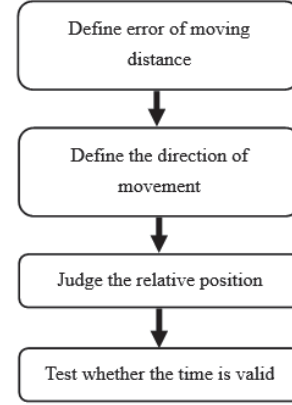


Fig. 4. Process of recognizing gestures

For right hand lifting up, the whole process is quite the same.

**Input:** Joint ID, maximum and minimum threshold value of time used to complete the whole gesture.

**Output:** The result of the recognition process

Initializing error in vertical direction, time range, and minimum distance.

start=0

for index from 1 to the length of skeleton linked list

//Determine which hand is moving

if(jointID==HandRight)

if list[index+1].Y - list[index].Y <0

start=index

//Determine the hand is moving vertically

if list[index].X - list[start].X >error in horizontal direction

start=index

//Determine the gesture is happening

if list[index].Y-list[start].Y>minimum distance and duration is in the range

return true

return false

Algorithm. 2. Recognition for right hand lifting up

#### IV. EXPERIMENTS

This paper has introduced a dynamic gesture recognition method based on Kinect. At present, there are three pre-defined basic gestures in standing postures: sliding left hand towards right, sliding right hand towards left and lifting up right hand. Experiments have been completed for testing the accuracy of recognizing these three gestures.

A system is designed for the training and testing processes. At the same time, the user-friendly interface of the system is designed with completely interactive function. Users can see the results directly from the interface. In addition, the rotation of a turntable represents the directions of dynamic gestures and displays the recognition results of input gestures.

When succeeding to recognize the specific gesture completed in a certain period of time, the system will sent out



Fig. 5. The interface of the system

the corresponding instructions to the targets. And the turntable is designed to show the recognition results. For instance, if the input gesture is recognized as left hand waving to the right, the turntable will rotate anticlockwise. If the gesture is right hand lifting up, the displaying result is a “STOP” sign (as the picture follows).

In the research of dynamic gesture recognition, the accuracy of recognition is one of the most important evaluation protocols. In order to get a convincing result, we have conducted a series of experiments to find out the appropriate value of parameters. For hands moving horizontally, we set minimum distance of hand movement as 0.35m and error in vertical direction as 0.12m. For right hand moving vertically, the minimum distance of hand movement is 0.45m and error in horizontal direction is 0.15m. The duration of all the gesture should be within 0.25s-0.75s. Each

TABLE I THE RESULT OF GESTURE RECOGNITION BASED ON KINECT

Gesture	Group	No. Tests	No. Correct Recognition	Accuracy	Avg Accuracy
Left hand waving through to the right	1	100	98	98%	97.67%
	2	100	98	98%	
	3	100	97	97%	
Right hand waving through to the left	1	100	95	95%	95.00%
	2	100	94	94%	
	3	100	96	96%	
Right hand lifting up	1	100	97	97%	97.33%
	2	100	98	98%	
	3	100	97	97%	

gesture is performed 300 times by three different people. And

we record the number of correct recognition and calculate its accuracy and average accuracy.

The experimental results prove that our system has a good performance for recognizing dynamic gestures. Our system defines recognizable gestures strictly. When the tested gesture is not in standard, this gesture may not be recognized. In addition, our system sets the time between two consecutive gesture recognition must be more than 0.5 seconds. If the gestures are made frequently, they may not be recognized correctly. In summary, three dynamic gestures' recognition accuracy is over 95%.

During the whole experiment, there are some reasons that conduce to such high rate.

- 1) Compared to traditional image recognition, skeleton data are more likely to recognize hands precisely. Because skeleton data use 3D information, which is less affected by environment and light, etc.
- 2) The predefined gestures are relatively easy to recognize. There are less differences among gestures performed by different people.
- 3) The parameters used in the experiment, such as time threshold, make the results convincing as well as avoiding interference of two continuous gestures.
- 4) Take errors into consideration, which give rise to high recognition rate.

Besides, high performance of Kinect, efficiency of the algorithm and standard test are all important factors in the experiment.

## V. CONCLUSIONS AND FUTURE WORK

In this paper, a dynamic gesture recognition method based on Kinect is introduced. Depending on the depth image and skeleton data of human bodies produced by Kinect, dynamic gestures can be captured and recognized efficiently. After determining which hand is moving, the method will give the result considering all the following factors: the direction and duration of movement. The experimental results show that our method is able to recognize the pre-defined dynamic gestures with a high recognition accuracy. While, the accuracy of the experiment will be influenced by the background, uncertainty of the gestures, mutual overlapping of humans' skeletons and many other factors. Our future work will be carried out to design more robust and efficient algorithms to recognize bimanual gestures. Besides, more gestures will be defined, which will make our system more practical. After developing a mature method, we plan to connect the system with intelligent appliances in houses and make it easier for the elderly to control these appliances.

## ACKNOWLEDGMENT

This paper is a funded by National College Student Innovation Experiment Program (201510561068), Guangdong Ministry of Education Foundation (No.2013B090500093), The Fundamental Research of Funds for The Central Universities (x2jsD2155000), Guangdong Province Science and Technology Plan Projects (2015A020219001), and Tianhe District science and technology project (201201YH038).



## REFERENCE

- [1] Binghai Ren, Yuanxin Guang et al. Vision-Based Recognition of Hand Gestures. Chinese Journal of Electronics 2000,28(2):118-121。
- [2] Youshu Hu, Survey of gesture recognition technology. China Science and technology information, 2005, 2:41-42。
- [3] ZhiQuan Feng, Yan Jiang. A survey of Hand Gesture of recognition. Journal of University of Jinnan, 2013, 27(4): 336-341。
- [4] S. S. Rautaray, A. Agrawal. Vision based hand gesture recognition for human computer interaction: a survey. Artificial Intelligence Review, 2012, DOI 10.1007/s10462-012-9356-9.
- [5] C. Charayaphayan, A. Marble. Image processing system for interpreting motion in American sign language. Journal of Biomedical Engineering, 1992, 14: 419-425.
- [6] T. Starner et al. Visual recognition of American Sign Language using hidden markov models. MIT Media Laboratory, MIT, USA, 1995.
- [7] G. Bradski, B. L. Yeo, M. M. Yeung. Gesture for video content navigation. Proceedings of Storage and Retrieval for Image and Video Database VII, 1999: 230-242.
- [8] M. J. Jones, J. M. Rehg. Statistical color model Models with Application to Skin Detection. IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'99), 1999, 1(1): 1274.
- [9] T. Starner, J. Weaver, A. Pentland. Real-time American Sign Language recognition using desk and wearable computer based video. IEEE transactions on pattern Analysis and Machine Intelligence, 1998, 20(12): 1371-1375