Real-Time Human-Robot Interaction for a Service Robot Based on 3D Human Activity Recognition and Human-mimicking Decision Mechanism

Kang Li, Jinting Wu, Xiaoguang Zhao and Min Tan
The State Key Laboratory of Management and Control for Complex Systems,
Institute of Automation, Chinese Academy of Sciences,
University of Chinese Academy of Sciences, Beijing, China
Email: likang2014@ia.ac.cn, wujinting2016@ia.ac.cn, xiaoguang.zhao@ia.ac.cn, min.tan@ia.ac.cn

Abstract-In this paper, we present a real-time Human-Robot Interaction (HRI) system for a service robot based on 3D human activity recognition and human-mimicking decision mechanism. A three-layer Long-Short-Term Memory (LSTM) network is trained to recognize human activity. The input of this network is the 3D skeleton joints information captured by Kinect V1 and the output of this network is the class label of human activity. Moreover, the human-mimicking decision mechanism is also fused into the online test process of human activity recognition, which allows the robot to instinctively decide whether to interrupt the current task according to task priority. The robot can make appropriate responses based on aforementioned analytical results. The system framework is realized on the Robot Operating System (ROS). The real-life activity interaction between our service robot and the user was conducted to demonstrate the effectiveness of developed HRI system.

I. INTRODUCTION

As robots move into our lives, the importance of enabling users to interact with them in a natural way increases [1]. Nonverbal communication has become the most natural manner between robots and users. Therefore, the robot with the capability of understanding human body language is a trend of development in HRI system.

In the past few years, gesture or activity recognition on RGB video processing has been widely studied and applied in various fields [2]-[6]. They attempted to learn discriminative features to enhance identification performance of human activities. However, hand-crafted features only have limited discriminative power, another drawback is that the 2D image only involves limited body movement information. With the rise of deep learning techniques [7], high-level semantic features of human activities can be extracted to improve the accuracy of human activity recognition. Especially the emergence of Recurrent Neural Networks (RNN) makes sequence learning easier [8]–[11]. On the other hand, 3D human activity information was also gradually applied to human activity recognition. Several researchers relied on stereo camera to obtain 3D position of human's head, hands or any other parts of body [12]–[14], Others encode the motion characteristics of an action by using depth maps from RGB-D camera [15]-

[17]. Using more rich 3D human activity information makes human activity recognition much easier and more accurate.

Unfortunately, the motivation of aforementioned approaches is that obtaining more accurate recognition results on a benchmark does not pay sufficient attention to the cost of time. And now there is no any one method can achieve 100% recognition accuracy, if the robot identification failed, corresponding wrong interactive task will go on until finished, which is time-consuming and power-hungry. Besides, if the robot is performing a task, a higher priority task is coming, the existing HRI system based on human's activity recognition will have to wait for the end of the current task, which looks very silly.

This work demonstrates progress towards a HRI system for a service robot capable of understanding common multimodal human activities in real time and determining whether to interrupt current interactive task cleverly. The 3D data of the skeleton joints consist of rich body movement information, which is collected from users via Microsoft Kinect and is used to train a three-layer Long-Short-Term Memory (LSTM) [7] network for user's activity recognition. Our previous work [18] elaborated the model and learning method. In this work, we described a interaction switch based on filter model of attention theory, simple and efficient gesture tells the service robot to start interaction. Combined with the interaction switch, human-mimicking decision mechanism of interrupting interactive tasks was also designed. Interruption can be triggered by the user when the interaction switch turns on, human-mimicking decision - "priority should be given to important urgent matters" will be executed. The overall HRI system was implemented into the Robot Operating System (ROS) and evaluated on a service robot. Results from humanrobot interaction experiments showed that aforementioned well-designed rules enabled the robot to become more smart.

The remainder of the paper is organized as follows. Section II describes the robot platform for HRI. Section III contains an overview of the human-robot interaction system framework and elaborates the details of method. Section IV provides the experimental results and analysis. Finally, Section V concludes this paper.

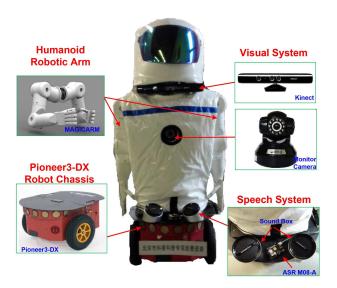


Fig. 1: Structure of service robot.

II. ROBOT PLATFORM AND CONFIGURATION

Our service robot adopts a distinctive modular concept to simplify mechanical design and installation. As shown in Fig. 1, the service robot consists of four principal components according to function: speech system, visual system, humanoid robotic arm, and Pioneer3-DX robot chassis.

- 1) The speech system includes an ASR M08-A module and two mini sound, which can identify speech command and broadcast the response, respectively.
- 2) The visual system includes a Microsoft Kinect and a remote monitor camera, which is used to observe and perceive the surrounding environment. Microsoft Kinect utilizes light coding to measure depth information. The measuring range can be up to approximately 3.5 m and the depth resolution can be up to 2 mm at a distance of 2 m. The field of view (FOV) is the vertebral rectangle, the horizontal view closes to 58°, and the resolution of the depth image is 640×480 pixels. The remote monitor camera is a common webcam and mainly responsible for monitoring.
- 3) In order to design friendly responses of the service robot when the user greets the robot, two humanoid robotic arms are installed on the service robot platform. The arms and hands have twelve and ten degrees of freedom, respectively. Generally, the two arms are 1.26 m long and weigh 15 kg. Each joint is a digital steering engine driven by pulse width modulation (PWM).
- 4) The Pioneer3-DX robot chassis is the main motion component, which has one follower wheel to maintain the balance of the robot and two driven wheels to control moving. In addition, the chassis is equipped with 16 sonar detectors to perceive obstacle of environment.

In this study, we only use part of the hardware comprising Microsoft Kinect, the robotic arms, and the Pioneer3-DX chassis for HRI.

III. HUMAN-ROBOT INTERACTION SYSTEM

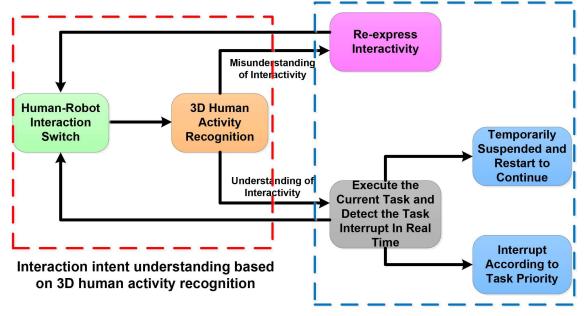
A. Overview of the framework of human-robot interaction system

Fig. 2 shows the framework of human-robot interaction system. The red dashed box for understanding the interactive intent of the service robot mainly includes two modules of human-robot interaction switch and 3D human activity recognition. Human-robot interaction switch uses the RGB-D camera to obtain the human skeleton point data stream, and opens the human-robot interaction according to well-designed natural gesture. Human activity recognition module trains a three-layer LSTM recurrent neural network and uses it to infer the meaning of human activity online. Blue dashed box for the task interruption mechanism mainly includes two strategies. One is the handling method when the interactive intent is misunderstood, the other is the coping approach of sudden task disruption during the execution of current correct task. In our method, when the intent of interactivity is misunderstood, combining with the designed human-robot interaction switch and 3D human activity recognition method, it is able to interrupt the wrong task execution and re-order the robot. When the intent of interactivity is correctly understood, the tasks will be executed and the disruptions will be detected in real time. Once the interruption occurs, robot will respond interruption according to priority of the task. Such humanmimicking decision mechanism is fused into HRI system, which realizes task suspension and switching, and enhances the intelligence and flexibility of interaction.

B. Human-robot interaction switch based on filter model of attention theory

1) The Filter Model of Cognitive Psychology Attention Theory: In cognitive psychology, the kernel of attention is the choice analysis of information. D.E.Broadbent believed that the information from the outside world is large, while the human nervous system has a limited central processing power. In order to avoid overloading the system, some sort of filter is required to adjust it and select the lesser and more critical information to move into the advanced analysis phase. Such information will be further processed and identified and stored, while other information will not be accepted. This is the theoretical filter model [19].

In the process of human-to-human interaction, we unconsciously notice each other's actions and filter out many less meaningful messages. In the process of human-robot interaction, in addition to the robot's ability to understand human's interactivity naturally and efficiently, it is also necessary for the robot to pay attention to human actions in real time during the mission, filter out useful information, such as people want to interrupt the current task of the intention. In our HRI system, the user can lift the left hand to attract the robot's attention. This simple and reliable skeleton-based interaction switch can make the robot filter out most of the worthless information, especially in multi-player scenarios.



Human-like intelligence decision mechanism based on attention model

Fig. 2: The framework of human-robot interaction system.



Fig. 3: Human skeleton point distribution.

2) The Implementation of Human-Robot Interaction Switch: Human skeletal point information refers to the spatial position coordinates of the vital joint point of the human body relative

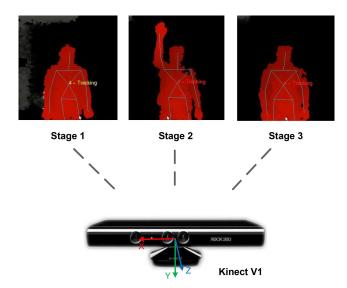


Fig. 4: The implementation process of human-robot interaction switch.

to the camera coordinate system collected by the RGB-D camera. As shown in Fig. 3, the human skeleton points collected by Microsoft Kinect mainly include the head, neck, torso, left wrist, left elbow, left shoulder, left hip, left knee, left ankle, right wrist, right elbow, right knee and right ankle total of 15 sets of data, and each set of data includes x, y, z coordinates of the three directions in the camera coordinate system. Therefore, the entire human skeleton point data can

be expressed as a set of 45-dimensional vectors:

$$F_n = [x_1, y_1, z_1, x_2, y_2, z_2, ..., x_{15}, y_{15}, z_{15}]$$
 (1)

where F_n represents the human skeleton point data in the n-th frame.

In the theory of attention in human's cognitive psychology, attention refers to the direction and concentration of a certain object by a mental activity. Following the characteristics of the notable object, we use the left hand gesture to turn on the interaction, which can stably and reliably attract the service robot's attention.

As shown in Fig. 4, the designed human-robot interaction switch is divided into three stages. In the first stage, the person stood normally and both arms fell naturally. At this point, both coordinate values of the left wrist and the right wrist of the human in the y direction are greater than the coordinate value of the human torso in the y direction. Here, the flag is set to false. In the second stage, the person raised his left arm. The coordinate value of the right wrist in the y direction is still greater than the coordinate value of the human torso in the y direction, but the coordinate value of the left wrist in the y direction is smaller than the coordinate value of the human torso in the y direction. Thus, the flag is set to true. In the third stage, the person put down the left arm. According to the coordinate value of the left wrist in the y direction is greater than the coordinate value of the human torso in the y direction and the flag is true, the robot would begin to interact with the person. The person can perform a body activity in 3.5 s, the data during this time would be recorded and fed into human activity recognition module.

C. Interactive intent inferring based on human activity recognition

Fig. 5 illustrates the basic flow of interactive intent inferring based on human activity recognition. First, turning on the human-robot interaction switch to collect the skeleton data of user's activity. Then, the skeleton data will be pre-processed by affine transformation to reduce the sensitivity from differences of user's location. Finally, we adopt LSTM networks to model human activity and determine user's interactive intent. The more details of three-layers LSTM model can be found in our previous work [18].

D. Motion control of robotic arm and chassis

After acquisition of user's interactive intent, the service robot responds correspondingly. Table. I shows eight types of interaction between the user and the robot. Here are two motion responses we designed for the robot: the motion of the robotic arm (Interaction 1, 3, 4 and 8) and the movement of the chassis (Interaction 2, 5, 6 and 7). For robotic arm motion control, the kinematics model is established using the Denavit-Hartenberg method [20]. Concrete parameters of the model can be found in our previous work [21]. For robotic chassis motion control, in order to stop the robot movement, the speed of chassis is set as 0. Moving forward and backward are both assigned fixed speed and mileage parameters,

circling is corresponding to spinning. These responses will be executed when the robot understands the corresponding human activities.

TABLE I: eight common types of interaction between user and robot

Human Activities (Interactive Intents)	Robot Responses
waving right hand (greet) stretching right hand (stop) saluting (salute) lifting right arm (lift right arm) waving forwards (go back) waving backwards (go ahead) drawing circle (perceive environment) waving arms around (march on the spot)	waving_right_hand stoping saluting lifting_right_arm moving_backwards moving_forwards circling waving_arms_around

E. Human-mimicking decision mechanism

In our HRI system, two cases of interactive task interruption are considered. One is to misunderstand the user's intent, the other is to understand the user's intent but need to interrupt current task. For the first case, the robot will carry out wrong task based on wrong interactive intent inferring. In the past HRI system, the user has to wait for the end of wrong task and restart the correct task. In our method, when the user observes that the robot misunderstand the user's intent, the current task can be suspended by the human-robot interaction switch. After the user lifts his left hand and attracts the robot's attention, the robot will stop the current task. When the user drops his left hand, the RGB-D camera will re-collect the user's body activity and re-infer the user's intent. This minor trick is consistent with human's smart decision mechanism. If others misunderstand the person's intent, the person will choose to re-express. For the second case, when the user's intent is correctly inferred, the robot will execute corresponding task. If you suddenly want to suspend the current task and re-start a higher priority task at this time, traditional interaction manner can not work. In our method, After the user lifts his left hand and draws attention of the robot, the robot will suspend the current task and record the breakpoint of the task. When the user put down his left hand, he can stop the robot using the stop command and handle his own thing, or uses the other command to control the robot. If stoping, the robot will forget the breakpoint and reset. If starting new interaction and the new command has higher priority, the new task will be opened and breakpoint will remember newest value. If starting new interaction and the new command has lower priority, the current task will be not interrupted and continued to execute. The overall algorithm has been described in Algorithm. 1.

IV. EXPERIMENTAL RESULTS AND ANALYSIS

A. Experimental Setup

We completed real-life human-robot interaction experiments in the hall including eight interactions between a human and the robot mentioned by Table. I. In our experiments, the parameters of human activity recognition model were set in

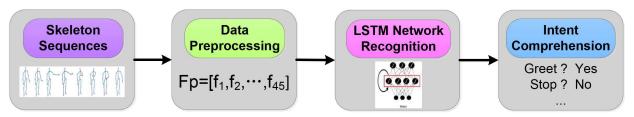


Fig. 5: The basic flow of interactive intent inferring based on human activity recognition.

2: opening human-robot interaction switch; 3: capturing skeleton data and feeding data into LSTM network to infer user's interactive intent; 4: if misunderstanding user's interactive intent then goto loop; 5: 6: else executing the current task and detect interruption in 7: real time: 8: if no interruption then 9: waiting the end of current task and **goto** loop; 10: else if new command is stoping then 11: forget the breakpoint and **goto** *loop*; 12: end if 13: 14: if new command has lower priority then interruption is invalid and current task will be 15: continued to execute;

if new command has higher priority then

interruption is valid and current task will be

end if

end if

suspended;

end if

16:

17:

18:

19:

20: **end if**

Algorithm 1 The procedure of human-robot interaction

[18]. In order to verify the effectiveness of adding humanmimicking decision mechanism for HRI system, we chose two simple task-based interactive mode, Interaction 5 and 7 shown in Table. I respectively. We defined that the task of moving back has higher priority than circling, and set the speed of the chassis as 0.2 m/s and the moving distance as 4 m, the center of the circle trajectory was set as the center of the chassis and the robot spins around 20 laps. First, the user executed drawing circle to tell the robot circling, the robot would rotate around itself. Then assuming that the task of moving backwards needed to do immediately, so the user lifted his left hand to interrupt current circling and performed waving forwards. After a while, we repeat aforementioned two activities alternately and observe the responses of the robot. Here we recorded the experimental results using screen recording mode in first view and video recording in third view.

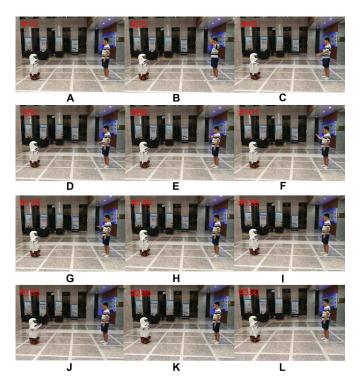


Fig. 6: Handwaving interaction between robot and user.

B. Experimental results and analysis

Here, we take handwaving interaction for example. Fig. 6 shows the results of the handwaving interaction. From A to C, the interaction switch was opened. From D to G, the user waved his right hand to greet the robot. From H to I, the robot analyzed the intention of user by calculating the output of LSTM network. From J to L, the robot responded user by waving its hand.

In the other experiment, in order to observe results more clearly, we chose a simplified robot platform without robotic arm because it does not use the robotic arm. Fig. 7 illustrates a series of results comprising of first view and third view. First, the user drew circle and the robot had a correct inferring. During the task, the user waved forwards, which has higher priority. The circling task is interrupted and the robot started moving back. After a while, the user tried to drew circle to interrupt the task of moving back, but failed because the circling has lower priority. When the user re-waved his hand forwards, the robot continued to move back from the

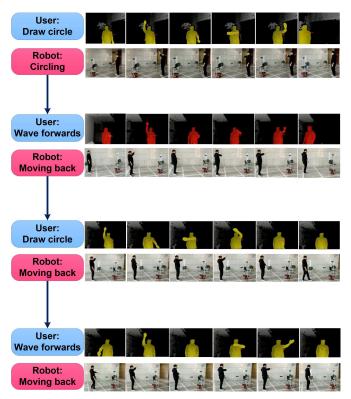


Fig. 7: The verification results of human-mimicking decision mechanism for the HRI system.

breakpoint. From the result, we can see that the robot has become more smart with the ability to determine whether to interrupt current task according to the priority of the task.

C. Discussion

The excellent performance of real-life experiments had illustrated the effectiveness of our HRI system. Based on the framework, an increasing number of interactive activities can be designed to enrich the interactive ability of the robot. However, our HRI system can only be used indoors because the skeleton data was captured by Microsoft Kinect V1. I believe that it is a challenging task that how to accurately capture human skeleton data in outdoors in real time.

V. CONCLUSION

In this paper, we described our progress towards a HRI system for a service robot capable of recognizing several daily user activities and making human-mimicking decision. A practical and natural interactive switch based on filter model of attention was designed for the HRI system. The human-mimicking decision – "priority should be given to important urgent matters" was also added into the HRI system. We presented results on several real-life interactive experiments, and demonstrated that our HRI system can realize excellent performance.

In the future, we will focus on the study of more intelligent interactive intent inferring for service robots. I believe that it is very interesting to make the robot more intelligent and smart.

ACKNOWLEDGMENT

This work is partially supported by the National Natural Science Foundation of China under Grants 61673378 and 61421004.

REFERENCES

- T. B. Sheridan. Human-robot interaction: Status and challenges. Human Factors: The Journal of the Human Factors and Ergonomics Society, 58(4):525–532, 2016.
- [2] Stefan Waldherr, Roseli Romero, and Sebastian Thrun. A gesture based interface for human-robot interaction. *Autonomous Robots*, 9(2):151– 173, 2000.
- [3] Ju Sun, Xiao Wu, Shuicheng Yan, Loong Fah Cheong, Tat Seng Chua, and Jintao Li. Hierarchical spatio-temporal context modeling for action recognition. In Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on, pages 2004–2011, 2009.
- [4] M. Sigalas, H. Baltzakis, and P. Trahanias. Gesture recognition based on arm tracking for human-robot interaction. In *Ieee/rsj International Conference on Intelligent Robots and Systems*, pages 5424–5429, 2010.
- [5] Quoc V Le, Will Y Zou, Serena Y Yeung, and Andrew Y Ng. Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis. 42(7):3361–3368, 2011.
- [6] Heng Wang, Alexander Klser, Cordelia Schmid, and Cheng Lin Liu. Dense trajectories and motion boundary descriptors for action recognition. *International Journal of Computer Vision*, 103(1):60–79, 2013.
- [7] Yann Lecun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. Nature, 521(7553):436–444, 2015.
- [8] Yong Du, Wei Wang, and Liang Wang. Hierarchical recurrent neural network for skeleton based action recognition. pages 1110–1118, 2015.
- [9] Wentao Zhu, Cuiling Lan, Junliang Xing, Wenjun Zeng, Yanghao Li, Li Shen, and Xiaohui Xie. Co-occurrence feature learning for skeleton based action recognition using regularized deep lstm networks. In *The* AAAI Conference on Artificial Intelligence, 2016.
- [10] Sijie Song, Cuiling Lan, Junliang Xing, Wenjun Zeng, and Jiaying Liu. An end-to-end spatio-temporal attention model for human action recognition from skeleton data. 2016.
- [11] Songyang Zhang, Xiaoming Liu, and Jun Xiao. On geometric features for skeleton-based action recognition using multilayer lstm networks. In IEEE Winter Conference on Applications of Computer Vision, 2017.
- [12] Rainer Stiefelhagen, Christian Fuegen, Petra Gieselmann, Hartwig Holzapfel, Kai Nickel, and Alex Waibel. Natural human-robot interaction using speech, gaze and gestures. 2004.
- [13] Nickel Kai and Rainer Stiefelhagen. Real-time person tracking and pointing gesture recognition for human-robot interaction. In *The* Workshop on Computer Vision in Human-Computer Interaction, pages 28–38, 2004.
- [14] Hee Deok Yang, A Yeon Park, and Seong Whan Lee. Gesture Spotting and Recognition for HumanRobot Interaction. IEEE Press, 2007.
- [15] Chen Chen, Kui Liu, and Nasser Kehtarnavaz. Real-time human action recognition based on depth motion maps. *Journal of Real-Time Image Processing*, 12(1):155–163, 2016.
- [16] Chen Chen, R Jafari, and N Kehtarnavaz. Action recognition from depth sequences using depth motion maps-based local binary patterns. In *IEEE Winter Conference on Applications of Computer Vision*, pages 1092–1099, 2015.
- [17] Chen Chen and Mengyuan Liu. 3d action recognition using multitemporal depth motion maps and fisher vector. In *International Joint Conference on Artificial Intelligence*, 2016.
- [18] Kang Li, Xiaoguang Zhao, Jiang Bian, and Min Tan. Sequential learning for multimodal 3d human activity recognition with long-short term memory. In *IEEE International Conference on Mechatronics and Automation*, pages 1556–1561, 2017.
- [19] D. E Broadbent. Perception and communication. Pergamon Press,, 1958.
- [20] R. S Hartenburg, J Denavit, and F Freudenstein. Kinematic synthesis of linkages. 34(30):no-no, 1965.
- [21] Kang Li, Ning An, Xiaoguang Zhao, Shiying Sun, and Min Tan. Body activity interaction for a service robot. In *IEEE International Conference* on Robotics and Biomimetics, pages 1942–1947, 2017.