

# Natural Teaching for Humanoid Robot via Human-in-the-loop Scene-motion Cross-modal Perception

**Abstract**—The paper aims to present a human-in-the-loop natural teaching paradigm based on scene-motion cross-modal perception, which facilitates the manipulation intelligence and robot teleoperation. The proposed natural teaching paradigm is used to telemanipulate a life-size humanoid robot in response to a complicated working scenario. First, a scene-motion cross-modal perception setup is built. A vision sensor is employed to project mission scenes onto virtual reality glasses for human-in-the-loop reactions. A motion capture system is established to retarget eye-body synergic movements to a skeletal model. Second, real-time data transfer is realized through publish-subscribe messaging mechanism in ROS (Robot Operating System). Next, joint angles are computed through a fast mapping algorithm and sent to a slave controller through a serial port. Finally, visualization terminals render it convenient to make comparisons between two motion systems. Experimentation in various industrial mission scenes, such as approaching flanges, shows the numerous advantages brought by natural teaching, including being real-time, high accuracy, repeatability and dexterity. The proposed paradigm realizes the natural cross-modal combination of perception information and enhances the working capacity and flexibility of industrial robots, paving a new way for effective robot teaching and autonomous learning.

**Index Terms**—Human-in-the-loop, Natural Teaching, Cross-modal Perception, Humanoid Robot, Motion Imitation

## I. INTRODUCTION

In recent years, demands for industrial robots with high intelligence have shown a tremendous growth in military, medicine, manufacturing and social life. Industrial robots are increasingly faced up with challenges of executing complicated tasks in unstructured environments, such as welding tracking on a curved surface, sorting and placing of scattered workpieces with surfaces of multiple types such as three-way valves and flanges. Many paradigms are adopted to improve the ability of robots to perform complex tasks based on data-driven methods [1]–[4]. However, with limited data, those data-driven methods alone tend to have a poor performance. Under such a circumstance, the combination of teaching and machine learning to cope with the lack of data has achieved good results [5]–[9]. As a direct way to endow industrial robots with humans knowledge, teaching renders the intelligence development of robots more than possible.

As a matter of fact, traditional teaching methods are faced with multiple difficulties. First, when teaching is performed for complicated motions with multiple DOF (degree of freedom), an expert is necessary for demonstration and the effect of teaching highly depends on his knowledge. The numerous frames of continuous movements will cause a sharp increase of the amount of teaching information, thus placing a heavy burden on the expert [10]–[12]. Second, in order to facilitate

robots to understand human teaching and to develop intelligence, behavior recognition and semantic classification are necessary [13], [14] while conventional demonstration methods often neglect the transmission of semantic information [15], [16]. Third, the ability to make decisions based on multiple sensory information is an important manifestation of human intelligence while conventional demonstration methods generally overlook or misunderstand the relation between different sensory information [17], [18].

Considering children's learning process, they observe the behavior of adults, and then reproduce it [19]. Such a process is always natural and highly effective because human beings share the same comprehension of scenes and the same behavioral language. Inspired by this, natural teaching is the key to overcoming obstacles to the exchange of teaching information between human and robots. Natural teaching is actually a branch of human-robot interaction (HRI) technology, representing a kind of teaching paradigm which is user-friendly and coordinates human and robot in scene comprehension. Aimed at completing tasks with specified human semantic information, natural teaching is an end-to-end and highly efficient method for interaction with surroundings or performing complicated movements. Moreover, training with such tasks is conducive to establish a deep understanding of potential implications from training data through subsequent intelligence algorithms, thus achieving a high level of intellectual development.

Scene-motion cross-modal perception constitutes a critical component of natural teaching. Inspired by role-play in esports, the demonstrator is provided with visual information to perceive the mission scenario of the robot and then implements various movements from a first-person view. The demonstrator's eye-body synergic movements are collected as motion information. Regarding to the teaching process, thanks to VR and HRI technology, the robot and the demonstrator can share the common visual and motion information during the whole process. In the aspect of teaching information, the robot achieves the cross-modal combination of scene and motion information with the assistance of human intelligence. The demonstrator can have an overall cognition of the surroundings from a first-person view, analyze complicated information and make movement decisions. The recording of intricate multi-DOF movements and the live video stream provide the robot with comprehensive scene-motion information so that the robot can be gradually endowed with the ability to repeat the same process. Further, the robot can even develop the capability of making autonomous decisions through such a natural teaching paradigm.

Employing humanoid robots as a platform to verify the natural teaching paradigm with scene-motion cross-modal perception can provide numerous advantages. First, since humanoid robots possess human-like structures and scales that have evolved for millions of years, the abundant DOF and the complex connections between links can characterize the control method of an industrial robot with an extremely complicated structure. Besides, the excellent mobility potential of humanoid robots renders it possible for them to be assigned with different tasks [20]. Second, humanoids can serve as a direct and natural platform for natural teaching. As they can completely reflect human motion, demonstrators can easily assess the difference between human motion and robot imitation during motion synchronization. Human can further consider the conversion of postures from human to robot and optimize the conversion rule against corresponding problems [21].

Herein we report a human-in-the-loop natural teaching paradigm with motion-scene cross-modal perception on a life-size humanoid robot. The robot is established based on InMoov, an open-sourced 3D printing humanoid robot. It posses a similar structure with human's and is equipped with 29 DOF, 22 of which are controlled in this system. The following is the natural teaching process. First, a vision sensor is employed to project the mission scene onto the VR glasses. Second, motion perception captures the motion of human with a set of wearable sensors and presents the collected motion data in BVH (BioVision Hierarchy) format. Then motion data are transmitted to an industrial PC (IPC) with Ubuntu running on it through TCP/IP and parsed according to BVH format. Next, the parsed euler angles are converted to corresponding joint angles through a fast mapping algorithm and encapsulated in a communication protocol. At last, IPC sends joint angles to the slave controller to control the robot. The whole system has paved a novel, real-time and accurate way for a natural teaching paradigm on humanoid robots.

This paper is organized as follows. In section 2, the scene-motion cross-modal perception system is introduced. Section 3 discusses the setup of the humanoid robot. Section 4 presents the realization of real-time motion imitation on the humanoid robot. Section 5 performs several experiments based on the proposed natural teaching paradigm. Finally, section 6 deals with the conclusion about our work.

## II. SCENE-MOTION CROSS-MODAL PERCEPTION

The framework of the cross-modal perception system is shown in Fig. 1. Scene perception makes it possible for the manipulator to perceive the complicated surroundings around the robot remotely, while motion perception passes back real-time human motions to the controller. The combination of scene and motion perceptions takes full advantage of humans intelligence because each movement in the loop is determined by human and reflected on the robot.

### A. Scene Perception

Scene perception is achieved through a remote video stream and multiple display terminals. Fig.2 shows its principle and

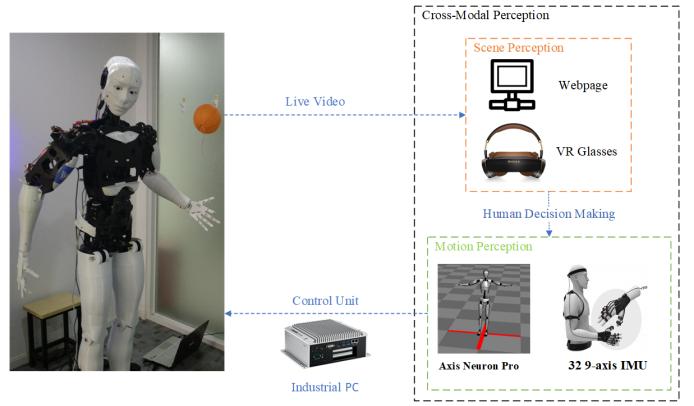


Fig. 1. Scene-motion Cross-modal Perception System

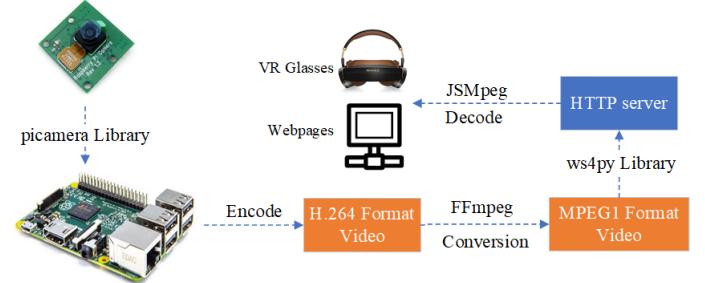


Fig. 2. Principle of Scene Perception

the video stream. The main function of scene perception is to stream the live video recorded by Raspberry Pi Camera to the manipulator to perceive the surroundings. Since the camera is installed in one eye of the robot, the manipulator wearing VR glasses can make decisions about movements from a first-person view. Besides, multiple display terminals make it possible for users to watch the same video stream on different electrical devices.

*1) Remote Video Stream:* Raspberry Pi is selected as the processing unit to drive the Pi camera for remote live video monitoring. The video obtained from raspberry is encoded in H.264 format, which is barely supported in browsers. Hence, FFmpeg (Fast Forward mpeg) is adopted to convert the H.264 format to the MPEG1 format. The video stream is then uploaded to a HTTP server through ws4py. To watch the video in the browser, the last step is decoding. JSMpeg is an excellent MPEG1 video and MP2 audio decoder defined in JavaScript. At most 30 fps video with resolution of  $1280 \times 960$  can be decoded by JSMpeg. Since JSMpeg is based on JavaScript, the video stream works in any modern browser (i.e. Firefox, Edge, Chrome, etc). Besides, the decoder has a low latency via WebSockets, thus achieving the real-time feature of our work.

*2) Multiple Display Terminals:* Display terminals include VR glasses and webpages. The type of VR glasses we adopt is Royole Moon, a combination of a headset, vari-focusing glasses and a control terminal. The operating system of Royole Moon is Moon OS which is developed based on Android.

Besides, it provides free access to external network, which means users can access the live video stream and perceive the mission scene at a distance. However, since the video capture is accomplished using one camera, all videos are 2D. Its inevitable that some necessary information will be lost by watching the screen alone. Therefore, its important to use VR glasses and some external assistances to improve the user experience. For webpage terminals, the principle is basically the same with VR glasses. Any devices which have installed a modern browser are accessible to the low-latency live video stream through a specified URL.

### B. Motion Perception

To capture motion information, several methods have been adopted. Abhay Bindal fixes IR sensors and accelerometer motion sensors to human legs and achieve real-time control of gaits on a biped humanoid robot [22]. Akif DURDU attaches potentiometers to human joints and then collect motion data [23]. Besides, vision sensing technology is also adopted. Several articles [24]–[26] utilize Kinect for gesture recognition and then perform similar actions on robots through different algorithms. Herein motion recording is achieved through wearable sensors. The motion capture system is composed of a motion sensor to capture real-time human motion and a human motion retargeting method.

1) *Motion Sensor*: A modular system composed of 32 9-axis sensors is adopted as the motion sensor. It is a set of wearable sensors designed by Noitom Technology Ltd to deliver motion capture technology. It contains 32 IMUs (Inertial Measurement Unit), each of which is composed of a 3-axis gyroscope, 3-axis accelerometer and 3-axis magnetometer. The static accuracy of each IMU is  $\pm 1$  degree for roll/pitch angle and  $\pm 2$  degree for yaw angle. The system is operated with Axis Neuron Pro (ANP) running on Windows OS for calibration and management. Besides, a skeleton model is visualized in ANP to reflect real-time human motion. Another important feature of ANP is to broadcast BVH data through TCP so that other programs can obtain and analyze these data using SDK provided by Noitom Technology Ltd.

2) *Human Motion Retargeting* : Motion retargeting is a classic problem which aims to retarget motion from one character to another while keeping styles of the original motion [27]. With this method, real-time human motion can be reflected on the skeletal model in ANP through BVH data. As a universal human feature animation file format usually adopted in skeletal animation models, it can store motion for a hierarchical skeleton, which means that motion of the child node is directly dependent on the motion of the parent one [28]. A normal BVH file will consist of several parts as follows.

- HIERARCHY signifies the beginning of skeleton definition.
- ROOT defines the root of the whole skeleton.
- OFFSET specifies the deviation of the child joint from its parent joint, which remains constant due to the unchanged lengths of human limbers.

```

HIERARCHY
ROOT Hips
{
    OFFSET 0.00 104.19 0.00
    CHANNELS 6 Xposition Yposition Zposition Yrotation Xrotation Zrotation
    JOINT RightUpLeg
    {
        OFFSET -11.50 0.00 0.00
        CHANNELS 6 Xposition Yposition Zposition Yrotation Xrotation Zrotation
        JOINT RightLeg
        {
            OFFSET 0.00 -48.00 0.00
            CHANNELS 6 Xposition Yposition Zposition Yrotation Xrotation Zrotation
            JOINT RightFoot
            {
                OFFSET 0.00 -48.00 0.00
                CHANNELS 6 Xposition Yposition Zposition Yrotation Xrotation Zrotation
                End Site
                {
                    OFFSET 0.00 -1.81 18.06
                }
            }
        }
    }
MOTION
Frames: 2
Frame Time: 0.04166667
-9.533684 4.447926 -0.566564 -7.757381 -1.735414 89.207932 9.763572

```

Fig. 3. An Example of BVH Format

- CHANNELS contains several parameters. The first parameter indicates the number of DOF. Usually only the root joint has both position data and rotation data. The rest ones only contain rotation data because positions of other joints can be obtained from OFFSET. Besides, these rotation data are Euler angles and the sequence of rotation hinges on the sequence mentioned in CHANNELS, i.e. the rotation is carried out in YXZ order in Fig. 3.
- End Site is only tagged in the definition of an end-effector and describes the lengths of bones through OFFSET.
- MOTION represents the beginning of another section which describes states of each joint at each moment.
- Frames stands for the current order of frames. Frame Time is the duration of each frame. The rest data are real-time states of each joint described sequentially in the HIERARCHY section. Hence, the number of these data is equal to the total number of channels defined in the HIERARCHY section.

We adopt BVH with no position channels. Hence three rotation values are obtained for each joint since position value keeps constant. Accordingly, we can figure out human gestures through these three rotation angles based on the assumption that wearable sensors are fixed with respect to human body.

### III. SETUP OF THE HUMANOID ROBOT

To realize real-time motion imitation on robots, a humanoid robot is set up since they possess human-like design and are able to mimic human motion [29]. However, due to the complicated structure of the robot and various constraints of conventional manufacturing methods, it is difficult to fulfill an elegant design of a dexterous humanoid robot. Fortunately, with the rapid advancement in 3D printing technology, 3D printing turns to be more cost-effective. Also, 3D printing element is also becoming more accurate, more complex and stronger. 3D-printed humanoid robots like InMoov, Flobi and

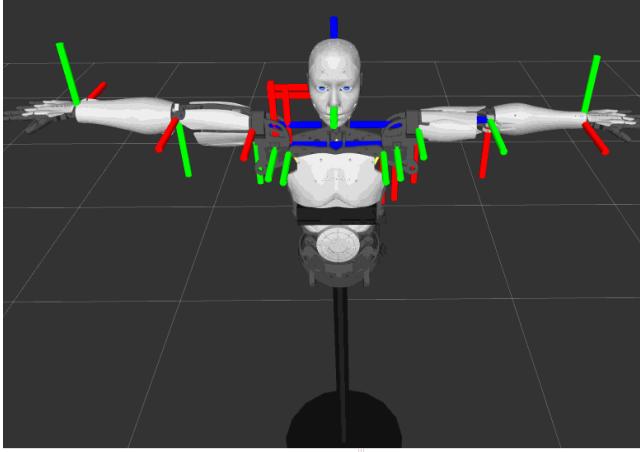


Fig. 4. DOF of Humanoid Robot (DOF of fingers are not displayed)

iCub have been created to serve as experiment platforms where research on HRI is conducted.

In this paper, a 3D-printed life-size humanoid robot is established based on InMoov initiated by Gael Langevin, a French sculptor in 2012 [30]. The whole structure as well as other necessary backgrounds have been illustrated in the previous work [31]. 22 out of 29 DOF are controlled during motion imitation, including 5 DOF for each hand, 4 for each arm, 3 for each shoulder and 2 for the neck, as shown in Fig. 4. As for control, the slave controller is composed of 4 small Arduino Nano control core boards, each of which can drive 6 servos with corresponding angles through PWM wave, and an Arduino Mega 2560 master board which communicates with the aforementioned nano nodes via 485 Hub based on the Modbus RTU control.

#### IV. REAL-TIME IMITATION OF HUMAN MOTION

The whole structure of the proposed method is shown in Fig. 5. First, the publish-subscribe messaging mechanism and the designed communication protocol ensures the security of data transfer. Second, the fast mapping algorithm converts BVH data into corresponding joint angles. Next, visualization terminals enables to make comparisons between different but simultaneous motion systems.

##### A. Data Transmission

During data transmission, scheduling protocols and quantization are required to prevent undesirable communication delays and packet dropouts [32], [33]. Herein the publish-subscribe messaging mechanism and a specified protocol are designed to realize the reliable data transmission. To be more specific, the publish-subscribe messaging mechanism allows nodes, which are executables after compilation, to publish messages or subscribe to a topic [34]. Topics are asynchronous and highly efficient. The whole data stream is mainly enabled through such a messaging mechanism, as shown in Fig. 6, where ellipses stand for nodes and squares represent topics.

- Socket\_node connects with the win32 console through TCP/IP and then advertises the topic, PN\_node/data.

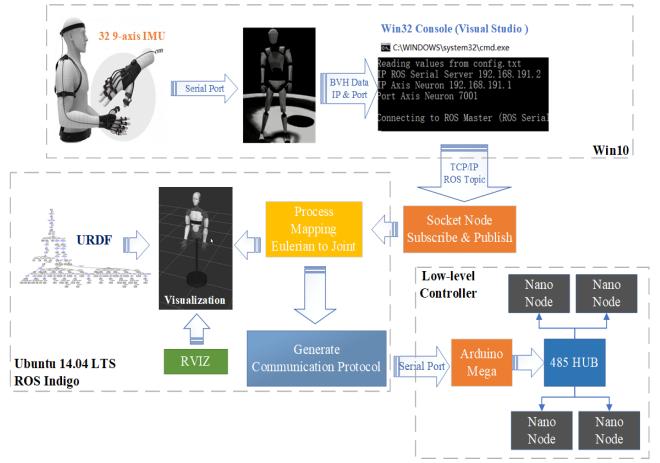


Fig. 5. Whole Structure of Proposed Method

- Mapping\_node subscribes to the previous topic and then converts BVH data to joint angles, which are then published to another topic called Joint\_angle.
- joint\_state\_publisher realizes the real-time simulation of robot model using the calculated joint angles.
- Serial\_node is responsible for the serial communication between the master and slave computers.

While topics have been successfully implemented in the data transfer process, reliable communication between the master and slave computers is still necessary to control the robot. Before transmission, all these data including a time stamp and joint angles are quantized to integers. The communication protocol contains 2 bits of time stamp data, 22 bits of position data corresponding to each joint, and 2 bits of CRC16 check code which are generated according to prior 27 bits to ensure the safety of data transfer.

##### B. Mapping Algorithm

Several methods have been adopted to achieve motion imitation. Y.Yuan utilizes a corrective wave variable method to address the force reflection control problem for motion imitation on a bilateral teleoperation system with time-varying delays [35]. M.Riley computes joint angles through a fast full-body inverse kinematics (IK) method [19]. The full-body IK problem is divided into many sub-problems to realize real-time imitation on a Sarcos humanoid robot with 30 DOF. J.Koenemann realizes complex whole-body motion imitation on a Nao humanoid based on the positions of end-effectors and center of mass [20]. By actively balancing the center of mass over the support polygon, the proposed approach enables the robot to stand on one foot as the demonstrator does. A.DURDU classifies the collected data with the assistance of ANN to perform movements on the robot [23]. Herein, a fast mapping algorithm is employed to realize the transformation.

To make the robot imitate human motion, the key point is to send corresponding joint angles computed from BVH data. BVH has provided us with three euler angles for each node,

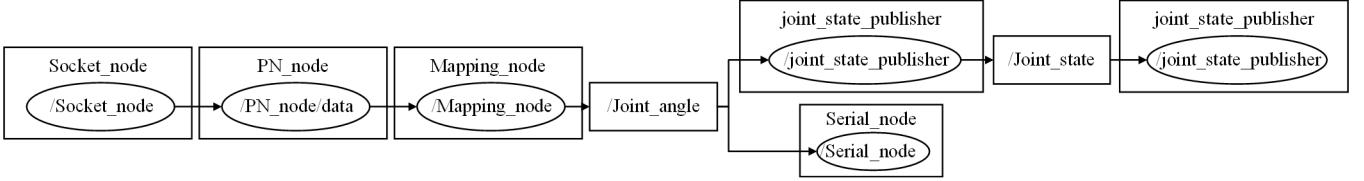


Fig. 6. Visualized Data Stream



Fig. 7. Designed Communication Protocol

enabling us to ascertain the rotation matrix between child and parent links. Denote euler angles with a rotation order of ZYX as  $\varphi, \theta, \psi$ , the rotation matrix of child frame with respect to parent frame is

$$R_{child}^{parent} = \begin{pmatrix} \cos\varphi & -\sin\varphi & 0 \\ \sin\varphi & \cos\varphi & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \cos\theta & 0 & \sin\theta \\ 0 & 1 & 0 \\ -\sin\theta & 0 & \cos\theta \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos\psi & -\sin\psi \\ 0 & \sin\psi & \cos\psi \end{pmatrix} \quad (1)$$

In this work, postures are mainly concerned for motion description. Hence, we consider human motion as a sequence of rotation matrices  $f_i$

$$f_i = \left\{ R_{LHand}^{LForearm}, R_{LForearm}^{LArm}, R_{LArm}^{Body}, R_{Head}^{Body}, R_{RHand}^{RForearm}, R_{RForearm}^{RArm}, R_{RArm}^{Body} \right\} \quad (2)$$

Thus, each posture is defined as a sequence of rotation matrices at time i, i.e.,  $R_{LHand}^{LForearm}$  stands for the rotation matrix between left hand and left forearm. Similarly, we can also define robot motion as another sequence. The goal is to eliminate the difference between each corresponding rotation matrix of human and robot as much as possible. Fig. 9 states the mapping problem. One one hand, human, with biological constraints, cannot have 3 rotational DOF at each joint and some of them are not completely independent. On the other hand, due to mechanical constraints, many joints of humanoid robot are also unable to rotate in three independent directions. Hence, each joint requires respective discussion for the mapping algorithm. Thanks to the structural symmetry, the algorithms for  $R_{LJoint2}^{LJoint1}$  and  $R_{RJoint1}^{RJoint2}$  share the same principle.

1) *Shoulder Joint*: The first case is the mapping between shoulders, which entails conversion from 3 human DOF to 3 robot DOF. Three rotational joints are installed on each shoulder part of InMoov and their axes of rotation can be approximately treated as perpendicular to each other. Denote the joint angles of 3 shoulder joints as respectively  $\alpha, \beta, \gamma$  and the rotation matrix of the arm link with respect to the body can be similarly expressed as

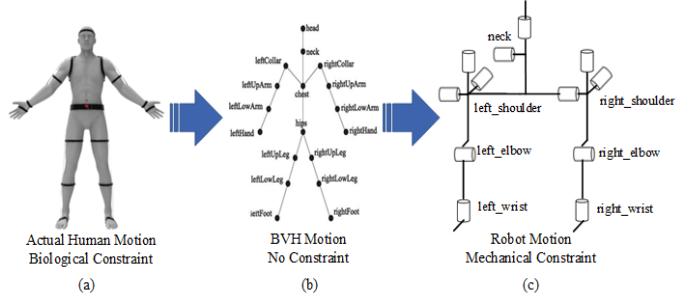


Fig. 8. Three Motion Systems with Different Constraints

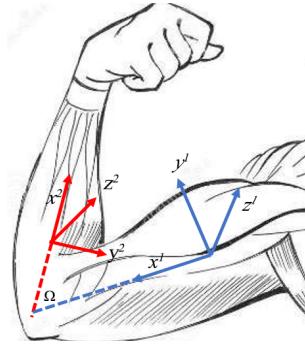


Fig. 9. Elbow Mapping

$$R_{Arm}^{Body} = \begin{pmatrix} \cos\alpha & -\sin\alpha & 0 \\ \sin\alpha & \cos\alpha & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \cos\beta & 0 & \sin\beta \\ 0 & 1 & 0 \\ -\sin\beta & 0 & \cos\beta \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos\gamma & -\sin\gamma \\ 0 & \sin\gamma & \cos\gamma \end{pmatrix} \quad (3)$$

With Equ.1 and Equ.3, we can derive a one-to-one correlation

$$\alpha = \varphi, \beta = \theta, \gamma = \psi \quad (4)$$

But there's still one thing that needs to be noticed. Since the mechanical structure of the robot has determined the rotation order of three joints between body and arm, the rotation order of euler angles in BVH should be set to be the same (in our case, ZYX).

2) *Elbow Joint*: The mapping between elbow joints entails the conversion from 2 human DOF to 1 robot DOF. Human elbows are able to bend and rotate while those of the robot can only bend. Then we need to compute the joint angle for bending, which is  $\Omega$ , as shown in Fig. 9. With the assumptions that sensors are fixed with respect to human body and the x-

direction is along the forearm link, we can derive the following equations.

$$\hat{\mathbf{x}_2}^2 = (1, 0, 0)^T \quad (5)$$

$$\hat{\mathbf{x}_2}^1 = R_2^1 \hat{\mathbf{x}_2}^2 = (\cos\varphi\cos\theta, \cos\varphi\sin\theta, -\sin\theta)^T \quad (6)$$

$$\langle \hat{\mathbf{x}_2}^1, \hat{\mathbf{x}_1}^1 \rangle = \arccos(\cos\varphi\cos\theta) \quad (7)$$

$$\Omega = \pi - \langle \hat{\mathbf{x}_2}^1, \hat{\mathbf{x}_1}^1 \rangle = \pi - \arccos(\cos\varphi\cos\theta) \quad (8)$$

$R_2^1$  stands for the rotation matrix of frame  $x_2y_2z_2$  with respect to  $x_1y_1z_1$ .  $\hat{\mathbf{x}_1}^1$  is a unit vector of  $\mathbf{x}_1$  in frame  $x_1y_1z_1$ .

3) *Neck Joint*: Mapping between neck joints requires the conversion from 3 human DOF( $\varphi, \theta, \psi$ ) to 2 robot DOF ( $\alpha, \beta$ ). Due to the mechanical constraints, rotation in one direction has to be abandoned. The solution to this case resembles that for the shoulder joint and be written as

$$\alpha = \varphi, \beta = \theta \quad (9)$$

### C. Visualization Terminals

On one hand, in the display of human motion collected from the motion sensor, ANP can visualize a skeletal model. Each joint is equipped with three DOF despite human's real physiological structures.

On the other hand, to visualize motion on the humanoid and to make simulation more convenient, another visualization scheme is provided with the assistance of ROS. A 3D visualization model is created in URDF (Unified Robot Description Format), a language based on XML and designed to describe the robot simulation model universally in ROS system, including the shape, size and colour, kinematic and dynamic characteristics of the model. The basic grammars are mentioned in [34]. However, the highly repetitive mechanical structure of InMoov makes it arduous to write a URDF manually. Hence we resort to a powerful tool called Xacro (XML Macros). Xacro is adopted to reuse the same structure for two different parts, i.e. left arms and right arms and to auto-generate a URDF file. Some fundamental grammars are shown in Table I. After importing \*.STL files with scale adjustment, the robot model can operate with the computed joint angles in RVIZ, a 3D visualization tool in ROS. These two visualization terminals are shown in Fig. 10.

## V. EXPERIMENT

This section designs several experiments to demonstrate the accuracy, repeatability and dexterity of the proposed natural teaching paradigm and discusses the experimental results.

### A. Accuracy

First, the accuracy of the control method is verified with several motion imitation experiments. Photos of various poses are taken using the motion imitation system, including various positions of two arms, face orientations and movements of fingers. The results can be examined in Fig. 11 -Fig.12. For these complicated gestures, the high degree of similarity between the wearer and the humanoid robot has demonstrated that the robot

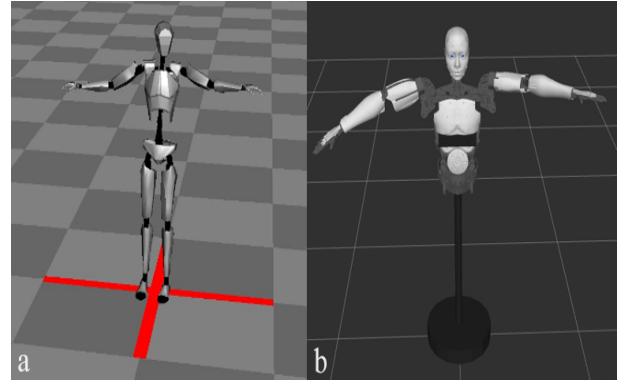


Fig. 10. Different Visualization Terminals for Different Motion Systems. (a) Skeletal Model in ANP. (b) Robot Model in RVIZ.



Fig. 11. Experiments of Different Gestures

has successfully followed the wearer's upper limber motion of the wearer, thus proving the feasibility and accuracy of our proposed method. Besides, the synchronous latency of less than 0.5 seconds validates the real-time performance.

To further illustrate accuracy, the second experiment is carried out to measure the angle errors for individual motions. Fig.13 shows the rotation directions and initial gestures. The angle errors for three motions are measured and plotted in Fig.14 and it follows  $Error = |\text{Angle}_{\text{Robot}} - \text{Angle}_{\text{Human}}|$ . The generally small errors are acceptable for natural teaching and further prove that the motion imitation system posses high accuracy.

However, there are still some limitations for motion imitation. The first one is the difference of structure between human and robot. Each of our arm has 7 DOF but the robot has only 5 and the rotational axes of their wrists are not the same. Besides, for some robot joints, the range of movement is limited due to its mechanical design. The second one is the mismatch between the skeletal model visualized through BVH data and the wearer's real motion. Revolution of each joint is achieved through skeletons in human body, while the wearable sensors can only remain fixed to the skin or clothes. The relative angular displacement between our skin and skeletons cause the measurement error. Other factors

TABLE I  
FUNDAMENTAL GRAMMARS OF XACRO

Command	Definition	Usage
Property	<code>&lt;xacro:property name="pi" value="3.14" /&gt;</code>	<code>&lt;... value =“ \${2*pi}”.../&gt;</code>
Argument	<code>&lt;xacro:arg name="use_gui" default="false"/&gt;</code>	<code>&lt;... use_gui:= true .../&gt;</code>
Macro	<code>&lt;xacro:macro name="arm" params="side"/&gt;</code>	<code>&lt;xacro:arm side="left"/&gt;</code>
Including	<code>&lt;xacro:include filename="other_file.xacro" /&gt;</code>	



Fig. 12. Comparison between Fingers

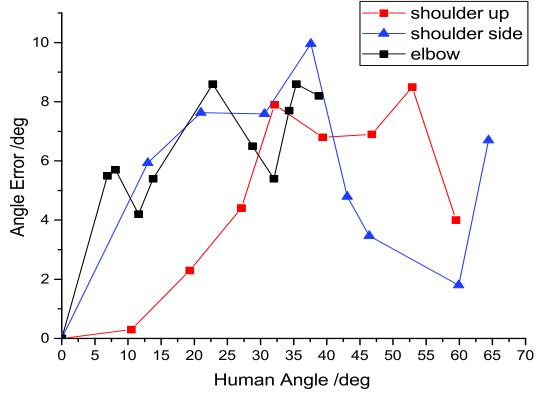


Fig. 14. Angle Error between Human and Robot

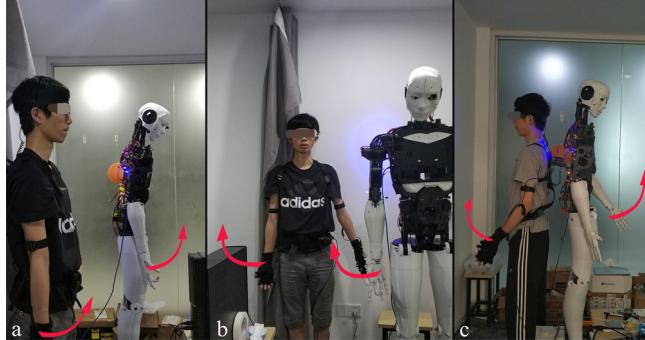


Fig. 13. Snapshots for Three Individual Motions. (a) Shoulder Up. (b) Shoulder Side. (c) Elbow.

include the accumulated drift errors and different positions of wearable sensors relative to human bodies. Nevertheless, there are still some possible solutions to these limitations. For example, sensors can be bound tightly to limbs in case of relative displacement between sensors and skin. Human motion can be confined to a certain range to achieve a higher accuracy. Besides, reasonable compensations for errors resulting from relative angular displacements between skins and skeletons can render the motion retargeting more reliable.

#### B. Repeatability

After demonstration is completed, the robot is expected to repeat the same motion, which means we high repeatability is desired. To verify the repeatability of the proposed natural teaching paradigm, an experiment where the robot is asked to approach the same point with its index finger is carried out. The distance errors are listed in Table. II. The average distance

error is 6.8mm, which is relatively small compared to its size. Aging of actuators, frictions in transmission mechanisms and instability of power supply may contribute to these gross angle errors.

TABLE II  
DISTANCE ERROR

Experiment No.	Distance Error		
	X <sup>a</sup> /mm	Y <sup>a</sup> /mm	R <sup>b</sup> /mm
1	0.0	0.0	0.0
2	2.6	-4.1	4.9
3	-14.3	2.2	14.5
4	-4.8	0.4	4.8
5	-8.4	-2.2	8.7
6	-10.4	-0.5	10.4
7	-13.8	-3.2	14.2
8	-7.3	2.3	7.7
9	-3.0	-1.2	3.2
10	0.0	0.0	0.0
Avg <sup>c</sup> /mm	-5.9	-0.6	6.8
Avg <sup>c</sup>  /mm	6.5	1.5	6.8

<sup>a</sup> The directions of X and Y are shown in Fig.15(a).

<sup>b</sup> R represents the distance error from the expected point and it follows  $R = \sqrt{X^2 + Y^2}$ .

<sup>c</sup> Avg is the abbreviation for average values.

To further illustrate the repeatability, another more compli-

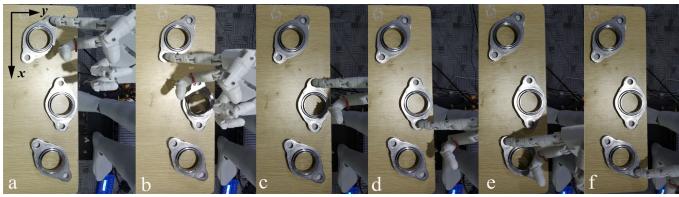


Fig. 15. Demonstration of Approaching Six holes

cated teaching experiment is carried out. Industrial robots are often required to repeat precise operations. Hence in the experiment, the demonstrator first teaches the robot to approach the left and right hole of 3 flanges in sequence and then the robot is asked to approach 6 holes continuously and automatically in each of the following experiments. The diameter of these holes is close to that of the fingertip. Fig. 15 shows the teaching process and Table. III shows the success rates. The experiment is carried out 5 times and failure mainly happens when the finger collides with the flange due to accumulated translational errors. The high success rate can validate the high repeatability and reliability of the proposed natural teaching paradigm. Once the demonstration is completed and desired movements or positions of the end-effector are also achieved, the robot is very likely to repeat the same motion using the recorded joint angles during the whole process.

TABLE III  
RESULTS OF REPEATING PROCESS

Experiment No.	a	b	c	d	e	f
1	S	S	S	S	S	S
2	S	S	S	S	S	S
3	S	S	S	S	S	F
4	S	S	S	F	S	S
5	S	S	F	S	S	S
Success Rate(%)	100	100	80	80	100	80

<sup>1</sup> a,b,c,d,e and f represent each hole in Fig.15

<sup>2</sup> F stands for failure and S stands for success

### C. Dexterity

With the proposed natural teaching paradigm, the humanoid robot is capable of accomplishing complicated movements. Two experiments are conducted to demonstrate the dexterity of natural teaching. First, an eye-body synergic experiment is performed via scene-motion cross-modal perception. In this experiment, the robot is faced with a complicated situation where flanges and other things heap up together on the table. As is the same with the aforementioned teaching experiment, the robot is asked to approach the inner circle of each flange with its index finger. The experimental results are shown in Fig. 16-Fig. 18.

The second one is the classical experiment of obstacle avoidance. As shown in Fig. 19, the end-effector needs to cross the obstacle first before it reaches the desired position

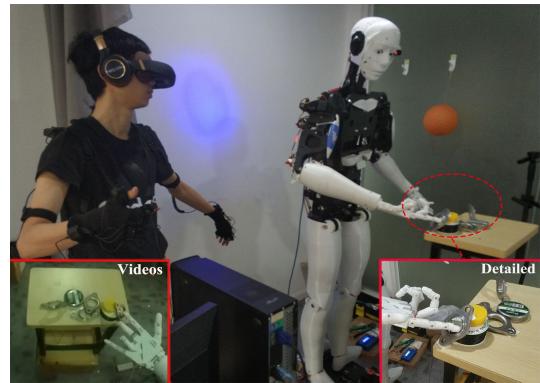


Fig. 16. Approaching an Inclined Flange



Fig. 17. Approaching a Vertical Flange

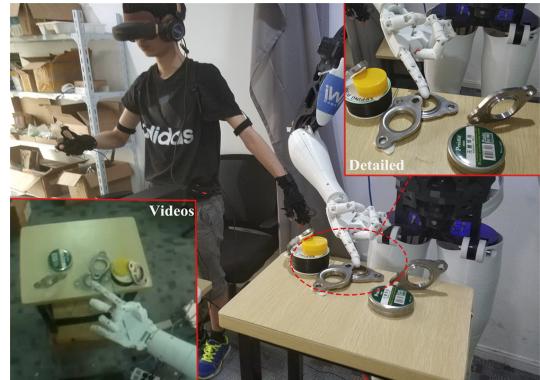


Fig. 18. Approaching an Occluded Flange

and orientation and the trajectory is generated by means of natural teaching. In such a complicated mission scene, motion planning always consumes a great amount of computation and time, while natural teaching can fully utilize humans perception and decision-making ability, making it more convenient and less time-consuming. The scene-motion cross-modal perception enables human to perceive the surroundings around the robot and the robot to reproduce humans motions. Hence, complicated missions can usually be accomplished with natural teaching at a minimum cost.

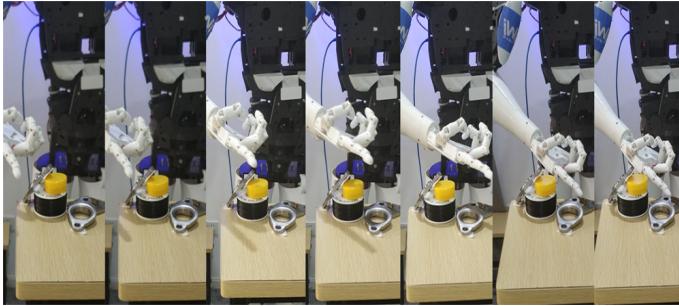


Fig. 19. Snapshots of End-effector Crossing oObstacles

## VI. CONCLUSIONS

This paper presents a novel natural teaching paradigm for a humanoid robot. A perception system composed of a vision sensor and a motion capture system realizes the cross-modal combination of scene and motion information. Through visualization terminals for different motion systems, a fast mapping algorithm and reliable data transfer methods, real-time motion imitation is accomplished. Several experiments are designed to validate the accuracy, repeatability and dexterity of the proposed natural teaching paradigm. Through natural teaching from a first-person view, human intelligence builds connections between scene information and movement policy, thus making it possible for robots to make autonomous decisions based on cross-modal perception. Future work will lay more emphasis on the development of the perception system to improve the user experience as well as the accuracy of motion imitation. Encouraged by Tri-Co Robot initiative [36], we hope this work will further contribute to the enhancement of industrial robot intelligence.

## REFERENCES

- [1] M. P. Deisenroth, G. Neumann, and J. Peters, *A Survey on Policy Search for Robotics*. Now Publishers Inc., 2013.
- [2] J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-driven grasp synthesis survey," *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 289–309, 2013.
- [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, and G. Ostrovski, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [4] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen, "Learning hand-eye coordination for robotic grasping with large-scale data collection," in *2016 International Symposium on Experimental Robotics*. Springer International Publishing, 2016, pp. 173–184.
- [5] N. Koenig and M. J. Matari, "Robot life-long task learning from human demonstrations: a bayesian approach," *Autonomous Robots*, vol. 41, no. 5, pp. 1–16, 2016.
- [6] C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine, "One-shot visual imitation learning via meta-learning," *ArXiv e-prints*, Sep 2017.
- [7] B. Michini, M. Cutler, and J. P. How, "Scalable reward learning from demonstration," in *IEEE International Conference on Robotics and Automation*, 2013, pp. 303–308.
- [8] Y. Kassahun, B. Yu, A. T. Tibebe, D. Stoyanov, S. Giannarou, J. H. Metzen, and E. V. Poorten, "Surgical robotics beyond enhanced dexterity instrumentation: a survey of machine learning techniques and their role in intelligent and autonomous surgical actions," *International Journal of Computer Assisted Radiology & Surgery*, vol. 11, no. 4, pp. 553–568, 2016.
- [9] S. Osentoski, "Remote robotic laboratories for learning from demonstration," *International Journal of Social Robotics*, vol. 4, no. 4, pp. 449–461, 2012.
- [10] D. Silver, J. A. Bagnell, and A. Stentz, "Active learning from demonstration for robust autonomous navigation," in *IEEE International Conference on Robotics and Automation*, 2012, pp. 200–207.
- [11] M. Laskey, J. Lee, C. Chuck, D. Gealy, W. Hsieh, F. T. Pokorny, A. D. Dragan, and K. Goldberg, "Robot grasping in clutter: Using a hierarchy of supervisors for learning from demonstrations," in *IEEE International Conference on Automation Science and Engineering*, 2016, pp. 827–834.
- [12] D. Koert, G. Maeda, R. Lioutikov, G. Neumann, and J. Peters, "Demonstration based trajectory optimization for generalizable robot motions," in *Ieee-Ras International Conference on Humanoid Robots*, 2017, pp. 515–522.
- [13] M. Wachter and T. Asfour, "Hierarchical segmentation of manipulation actions based on object relations and motion characteristics," in *International Conference on Advanced Robotics*, 2015, pp. 549–556.
- [14] G. H. Lim, "Two-step learning about normal and exceptional human behaviors incorporating patterns and knowledge," in *IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems*, 2017, pp. 162–167.
- [15] L. Rozo, S. Calinon, D. G. Caldwell, P. Jimnez, and C. Torras, "Learning physical collaborative robot behaviors from human demonstrations," *IEEE Transactions on Robotics*, vol. 32, no. 3, pp. 513–527, 2016.
- [16] M. Alibeigi, M. N. Ahmadabadi, and B. N. Araabi, "A fast, robust, and incremental model for learning high-level concepts from human motions by imitation," *IEEE Transactions on Robotics*, vol. PP, no. 99, pp. 1–16, 2017.
- [17] A. Lim and H. G. Okuno, "The mei robot: Towards using motherese to develop multimodal emotional intelligence," *IEEE Transactions on Autonomous Mental Development*, vol. 6, no. 2, pp. 126–138, 2014.
- [18] K. Noda, H. Arie, Y. Suga, and T. Ogata, "Multimodal integration learning of robot behavior using deep neural networks," *Robotics & Autonomous Systems*, vol. 62, no. 6, pp. 721–736, 2014.
- [19] M. Riley, A. Ude, K. Wade, and C. G. Atkeson, "Enabling real-time full-body imitation: a natural way of transferring human movement to humanoids," in *IEEE International Conference on Robotics and Automation, 2003. Proceedings. ICRA*, 2003, pp. 2368–2374 vol.2.
- [20] J. Koenemann, F. Burget, and M. Bennewitz, "Real-time imitation of human whole-body motions by humanoids," in *IEEE International Conference on Robotics and Automation*, 2014, pp. 2806–2812.
- [21] B. D. Argall, S. Chernova, M. Velso, and B. Browning, "A survey of robot learning from demonstration," *Robotics & Autonomous Systems*, vol. 57, no. 5, pp. 469–483, 2009.
- [22] S. Gobee, M. Muller, V. Durairajah, and R. Kassoo, "Humanoid robot upper limb control using microsoft kinect," in *2017 International Conference on Robotics, Automation and Sciences (ICORAS)*, Nov 2017, pp. 1–5.
- [23] A. Durdu, H. Cetin, and H. Komur, "Robot imitation of human arm via artificial neural network," in *International Conference on Mechatronics - Mechatronika*, 2015, pp. 370–374.
- [24] E. Yavan and A. Uar, "Gesture imitation and recognition using kinect sensor and extreme learning machines," *Measurement*, vol. 94, pp. 852–861, 2016.
- [25] I. J. Ding, C. W. Chang, and C. J. He, "A kinect-based gesture command control method for human action imitations of humanoid robots," in *International Conference on Fuzzy Theory and ITS Applications*, 2014, pp. 208–211.
- [26] A. Bindal, A. Kumar, H. Sharma, and W. K. Kumar, "Design and implementation of a shadow bot for mimicking the basic motion of a human leg," in *International Conference on Recent Developments in Control, Automation and Power Engineering*, 2015, pp. 361–366.
- [27] X. Meng, J. Pan, and H. Qin, "Motion capture and retargeting of fish by monocular camera," in *International Conference on Cyberworlds*, 2017, pp. 80–87.
- [28] H. Dai, B. Cai, J. Song, and D. Zhang, "Skeletal animation based on bvh motion data," in *2010 2nd International Conference on Information Engineering and Computer Science*, Dec 2010, pp. 1–4.
- [29] N. E. N. Rodriguez, G. Carbone, and M. Ceccarelli, "Antropomorphic design and operation of a new low-cost humanoid robot," in *Ieee/ras-Embs International Conference on Biomedical Robotics and Biomechatronics*, 2006, pp. 933–938.
- [30] G. Langevin, "Inmoov," <http://www.inmoov.fr/project>, 2014.

- [31] L. Gong, C. Gong, Z. Ma, L. Zhao, Z. Wang, X. Li, X. Jing, H. Yang, and C. Liu, "Real-time human-in-the-loop remote control for a life-size traffic police robot with multiple augmented reality aided display terminals," in *2017 2nd International Conference on Advanced Robotics and Mechatronics (ICARM)*, Aug 2017, pp. 420–425.
- [32] Y. Yuan, Z. Wang, P. Zhang, and H. Liu, "Near-optimal resilient control strategy design for state-saturated networked systems under stochastic communication protocol," *IEEE Transactions on Cybernetics*, pp. 1–13, 2018.
- [33] K. Liu, E. Fridman, K. H. Johansson, and Y. Xia, "Quantized control under round-robin communication protocol," *IEEE Transactions on Industrial Electronics*, vol. 63, no. 7, pp. 4461–4471, July 2016.
- [34] Z. Wang, L. Gong, Q. Chen, Y. Li, C. Liu, and Y. Huang, *Rapid Developing the Simulation and Control Systems for a Multifunctional Autonomous Agricultural Robot with ROS*. Springer International Publishing, 2016.
- [35] Y. Yuan, Y. Wang, and L. Guo, "Force reflecting control for bilateral teleoperation system under time-varying delays," *IEEE Transactions on Industrial Informatics*, pp. 1–1, 2018.
- [36] H. Ding, X. Yang, N. Zheng, M. Li, Y. Lai, and H. Wu, "Tri-co robot: a chinese robotic research initiative for enhanced robot interaction capabilities," *National Science Review*, p. nwx148, 2017. [Online]. Available: <http://dx.doi.org/10.1093/nsr/nwx148>