

Robot Gesture Generation from Environmental Sounds Using Inter-modality Mapping

Yuya Hattori* Hideki Kozima** Kazunori Komatani* Tetsuya Ogata* Hiroshi G. Okuno*

*Graduate School of Informatics,
 Kyoto University, Japan

{yuya, komatani, ogata, okuno}@kuis.kyoto-u.ac.jp

**National Institute of Information and
 Communication Technology, Japan

xkoizima@nict.go.jp

Abstract

We propose a motion generation model in which robots presume the sound source of an environmental sound and imitate its motion. Sharing environmental sounds between humans and robots enables them to share environmental information. It is difficult to transmit environmental sounds in human-robot communications. We approached this problem by focusing on the iconic gestures. Concretely, robots presume the motion of the sound source object and map it to the robot motion. This method enabled robots to imitate the motion of the sound source using their bodies.

1. Introduction

Based on advances in information technologies, the widespread use of robots is expected. Robots interacting with humans in the real world are supposed to do so using diverse modalities as humans do. In particular, we focus on various kinds of non-verbal sounds in our surroundings such as a door-opening sound and the cry of an animal, called “environmental sounds.” Since environmental sounds are very important clues to understanding the surroundings, it is very useful to share the information of the environmental sounds with humans and robots. Ishihara et al. developed a sound-to-onomatopoeia translation method for such interactions (Ishihara et al., 2004), for example. In this paper, we focus on gestures in order to interact about environmental sounds.

2. Motion Generation using Inter-modality Mapping

2.1 Definition of Inter-modality Mapping

Humans ordinarily perceive events in the real world as the stimuli of multiple modalities such as vision and audition. Accordingly humans can express the stimuli in each modality. In the real environment, information from all modalities is not obtained properly. Optical information can have occlusions, for example. In such cases, humans can complement losses from properly obtained information. We defined **inter-modality mapping** as mapping from information of properly obtained modalities to the information of modalities not obtained.

In this paper, we focus on mapping from input sounds to motions. It is because visual and auditory information are remarkably important in humans’ communications. It has been reported that children often use onomatopoeia and a gesture simultaneously and that the linkup of them is important for the development of multi-modality interactions (Werner and Kaplan, 1963).

2.2 Iconic Gesture Generation

We aim to generate motions expressing environmental sounds by imitating the motions of the sound source objects. It is because the kinds of environmental sounds are closely related to the motion of the sound source. It is known as **iconic gestures from an observer viewpoint** to imitate the motion of objects (McNeill, 1992). Iconic gesture means imitating concrete circumstances or events using one’s body.

In our model, robots memorize the motion of the sound source, which is captured by their camera, when they listen to a sound with looking the sound source. After learning the correspondences between the sound and the motion, they can imitate the motion of the sound source when they listen to a sound without looking at the sound source. Namely, they can perform iconic gestures.

3. System Implementation

3.1 Tasks and System Overview

In order to make iconic gestures, the system learns connections between the sound and the motion when a sound occurs. In interaction phase, the system generates a motion when a sound is input.

We use the robot “Keepon” for implementation and experiments in this study. It is a creature-like robot developed at NICT mainly for communicative experiments with infants. Its body is approximately 12 cm high.

3.2 Learning Process

If object velocity, or the norm of the optical flow vector, is higher than the threshold when a sound is input, the system interprets the motion of the object as the reason why the sound occurred. The system memorizes the pair of the spectrogram of the separated sound and the sequence of the optical flow vectors into the mapping memory. This process is shown in Fig. 1. The detailed process in each module is explained in the following paragraphs.

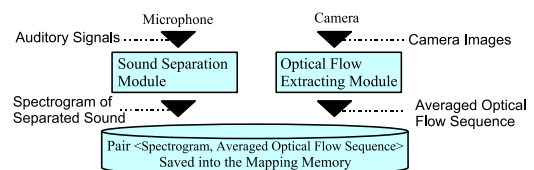


Figure 1: Motion and Sound Learning

Optical Flow Extraction Optical flows are constantly extracted from camera images. We adopted the block matching method for optical flow. Assuming that the camera captures only a sound source object, the system averages the flows of all of the blocks

Sound Separation In most physical phenomena, one sound is equivalent to one peak in the power envelope of an input audio signal. Therefore, each event is separated by extracting each peak. We adopted the separation method developed by Ishihara et al. (Ishihara et al., 2004). Individual sound separation is done by extracting each local maximum of power envelope which is not regarded as a part of another local maximum.

$$\mathbf{T}_i = [\min_t \{ t \mid V_s < |\mathbf{F}(t)|, t \leq \min \mathbf{T}_s \}, \\ \max \{ t \mid V_s < |\mathbf{F}(t)|, \max \mathbf{T}_s \leq t \}]$$

3.3 Sound-to-Motion Translation

```

graph TD
    Microphone -- Auditory Signals --> SSMS[Sound Separation Module]
    SSMS -- "Spectrogram of Separated Sound" --> ADCM[Auditory Distance Computation Module]
    ADCM <--> MM[(Mapping Memory)]
    ADCM -- "Optical Flow Sequence of the Sound Giving Minimum Distance" --> MGM[Motion Generation Module]
    MGM -- "Reproducing Optical Flow Sequence" --> Output
  
```

Auditory Distance Computation Auditory distances are computed by dynamic time warping (DTW) with mel filter bank output. DTW has been reported as one of the most effective methods in recognition of environmental sounds (Cowling and Sitte, 2003). We adopt mel filter bank output, which is made to adapt to human perception, as a feature of each frame in DTW.

$$\mathbf{X}(t') = C \sum_{s=1}^{t'} \mathbf{F}_i(s)$$

can reproduce the generated motion, must have been selected from all of the degrees of freedom of the robot beforehand.

The motion generation system was implemented on Keepon. We made Keepon learn 4 kinds of sounds, such as the sound of rubbing uneven plastic pieces together and the sound of striking metal boards vertically. After the learning, we made Keepon generate the imitation motion by giving a sound without looking at the sound source. Based on the rubbing sound, Keepon made a right rotation, then a left rotation as shown in Fig. 3 (a). According to the striking sound, Keepon executed an anteflexion then a retroflexion as Fig. 3 (b). The coordinates of these motions are shown as Fig. 4. A video of these experiments can be viewed at <http://winnie.kuis.kyoto-u.ac.jp/members/yuya/demo.avi>.



5. Conclusion

Acknowledgments This study was partially supported by the Grant-in-Aid for Scientific Research from the Japanese Ministry of Education, Science, Sports and Culture, and by the JPSP 21st Century COE Program.

Cowling, M. and Sitte, R. (2003). Comparison of techniques for environmental sound recognition. *Pattern Recogn. Lett.*, 24(15):2895–2907.

Ishihara, K., Nakatani, T., Ogata, T., and Okuno, H. G. (2004). Automatic sound-imitation word recognition from environmental sounds focusing on ambiguity problem in determining phonemes. *PRICAI-2004*.

McNeill, D. (1992). *HAND AND MIND: What Gestures Reveal about Thought*. Univ. of Chicago Pr, Chicago.

Ramachandran, V. S. and Hubbard, E. M. (2001). Synaesthesia — a window into perception, thought and language. *Journal of Consciousness Studies*, 8(12):3–34.

Werner, H. and Kaplan, B. (1963). *Symbol Formation: An Organismic-Developmental Approach to the Psychology of Language*. John Wiley and Sons, New York.