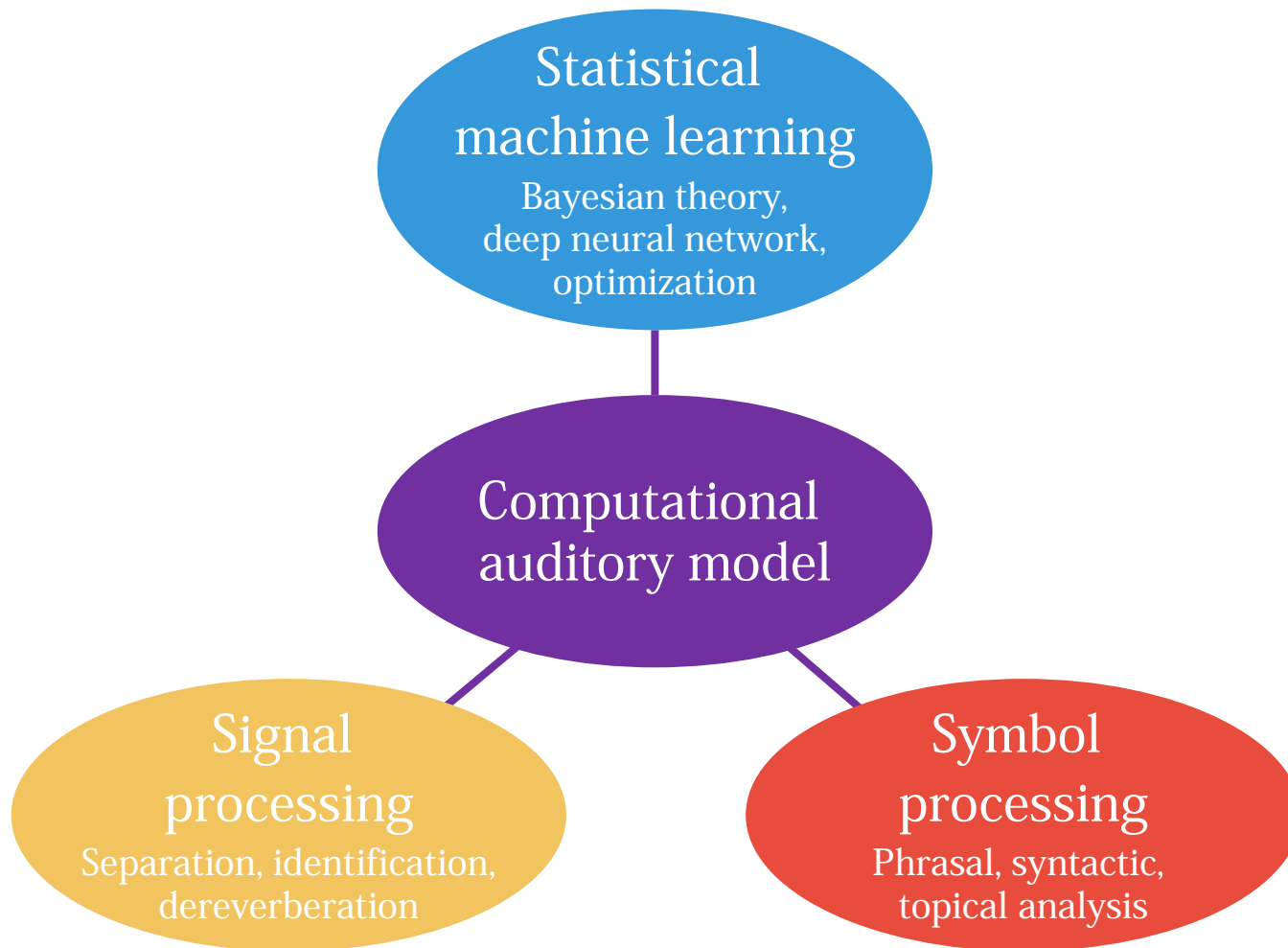


Audio Media Processing

**Graduate School of Informatics
Kyoto University**

Kazuyoshi Yoshii
yoshii@kuis.kyoto-u.ac.jp



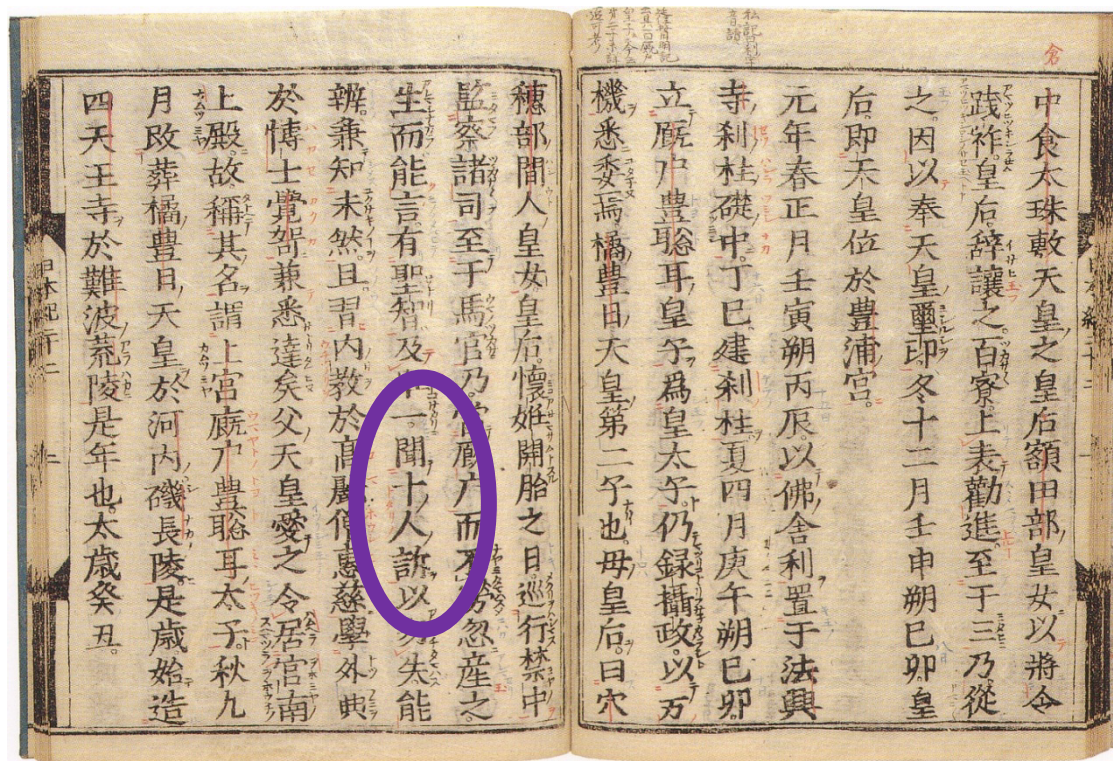
- Listen to “speech”
 - Prince-Shotoku robot
 - ♦ Simultaneous speech recognition
 - Microphone-array processing
 - ♦ Sound source separation and dereverberation
- Listen to “music”
 - Music understanding and performance
 - ♦ Sound source separation and music transcription
 - ♦ Co-playing and accompaniment
- Listen to “environmental sounds”
 - Object detection in a disaster environment
 - Analysis of frog calling

Listen to Speech

Shotoku-Taishi (Prince Shotoku)

5

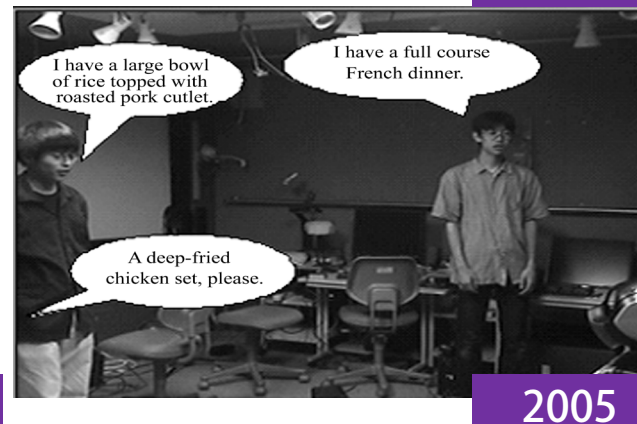
- Legendary person who can recognize simultaneous utterances of ten persons



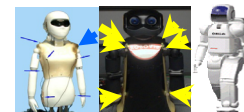
Simultaneous Speech Recognition

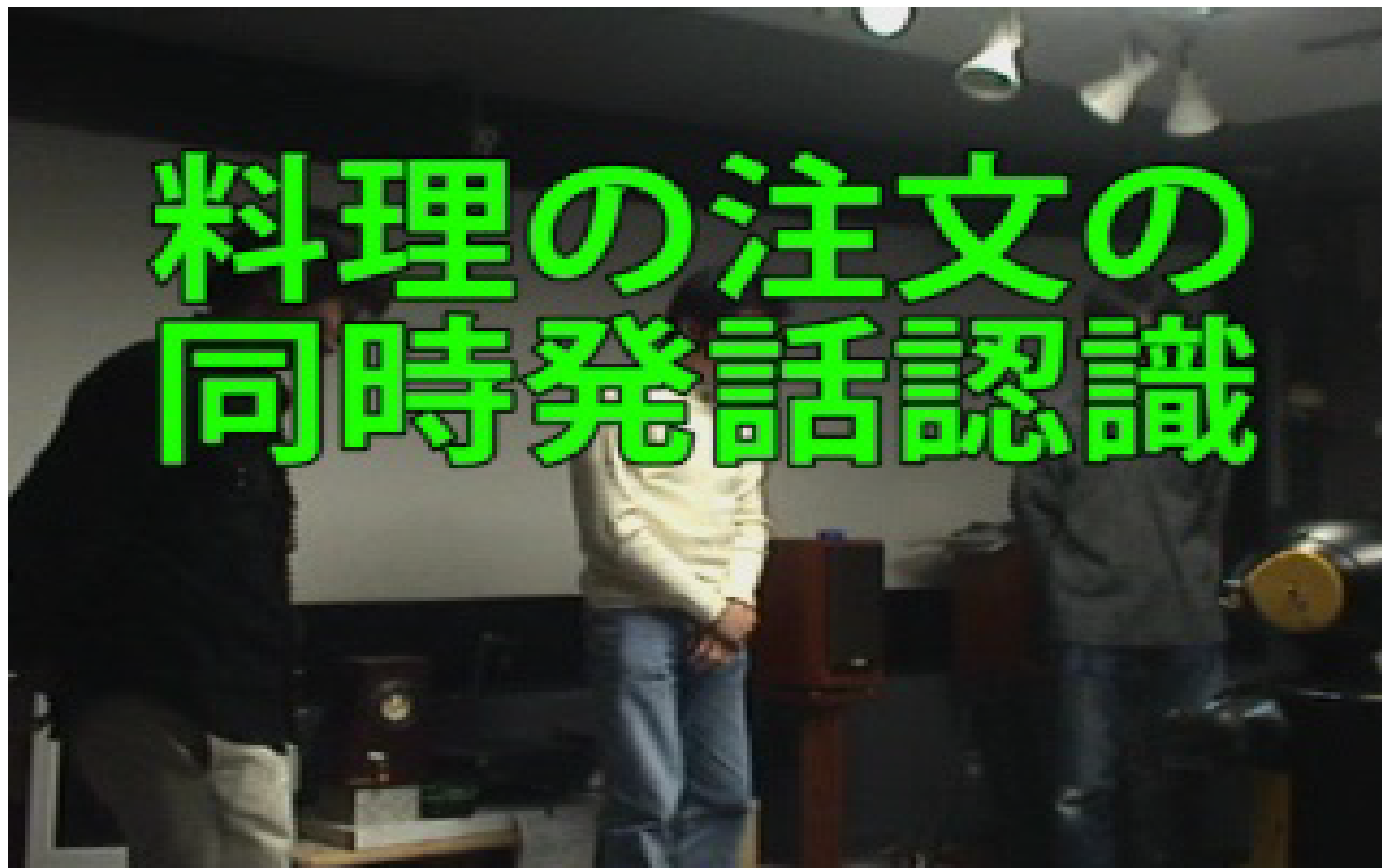
6

Closeness between speakers



- Isolated word → Continuous speech
- Evaluation using three robots in a large room





- Can many-ears robots go beyond humans?

Simultaneous Speech Recognition

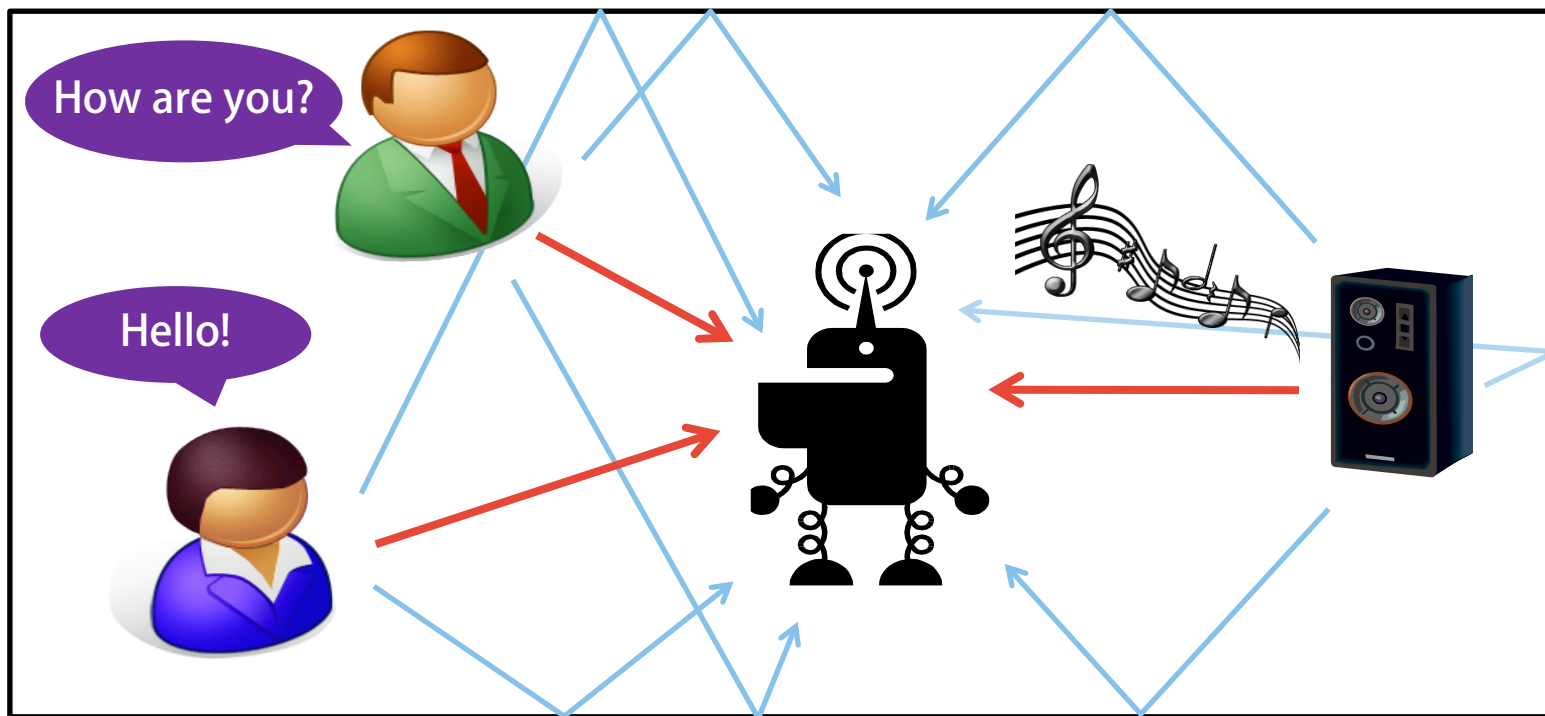
～ Meal Order Taking ～

16ch microphone-array processing • Sound directions are given

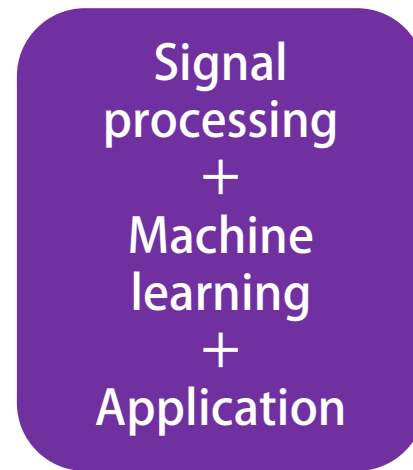
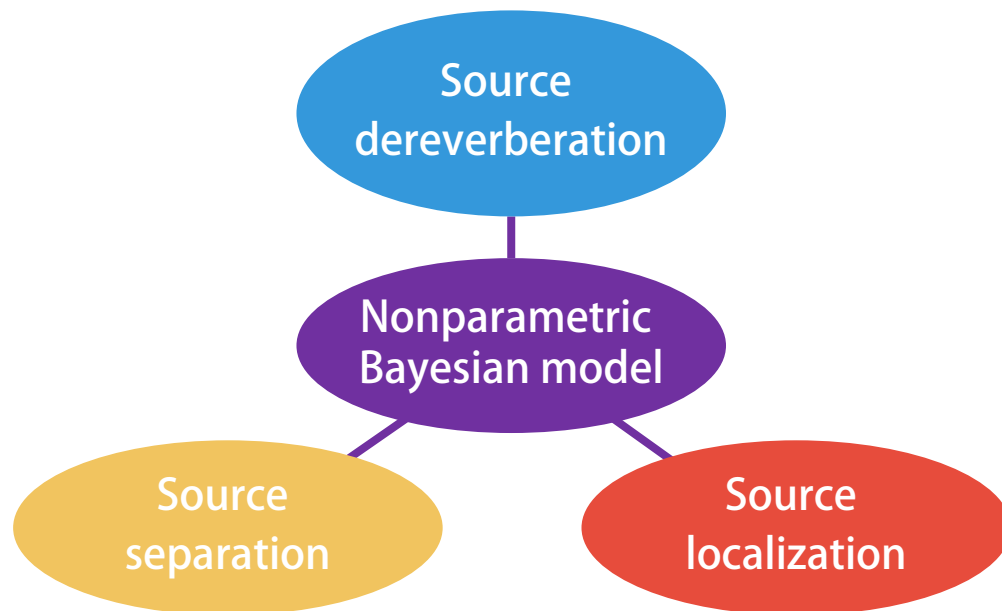
Microphone Array Processing

9

- Separate mixture signals into unknown number of sound sources with unknown reverberation time



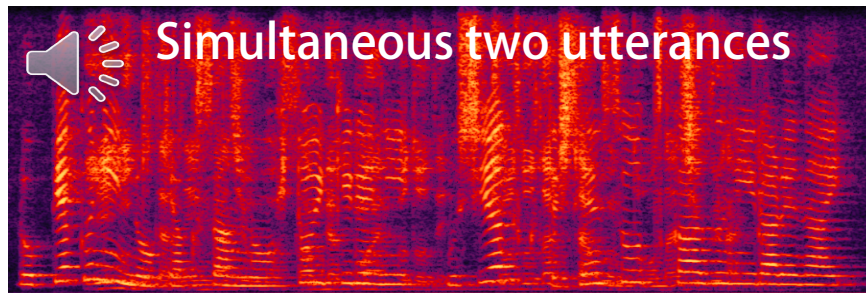
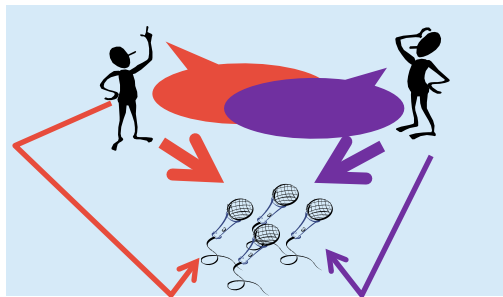
- A modern principled approach to the conventional egg-and-chicken problem
 - C.f. Errors are propagated in a cascade framework
(**localization** → **separation** → **dereverberation**)



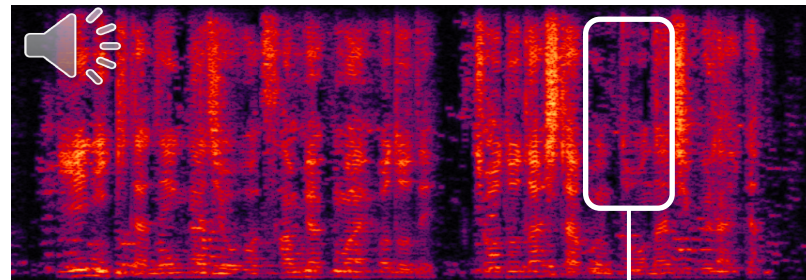
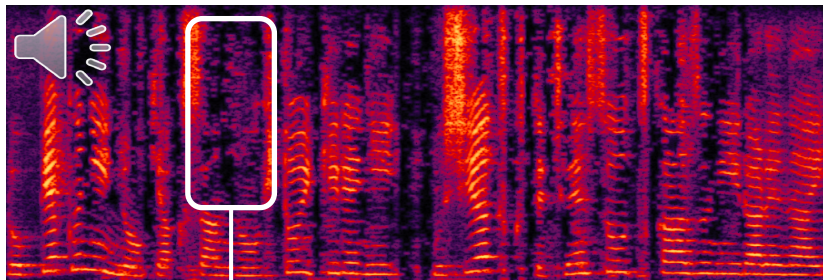
Hot topic!

Localization + Separation + Dereverberation

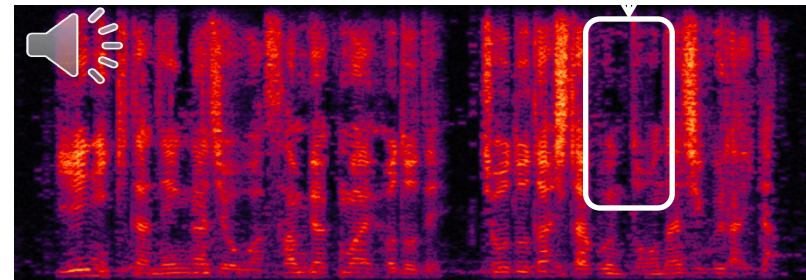
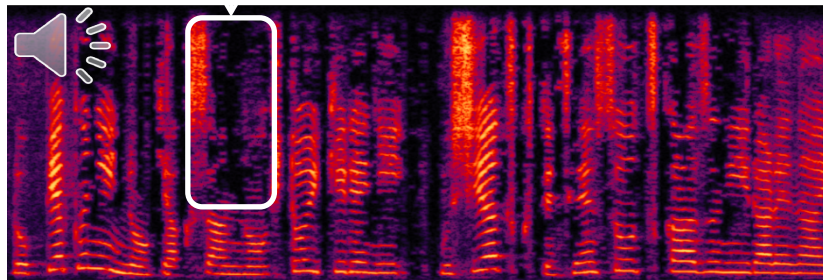
11



Without dereverberation



With dereverberation



Application to Real Environment

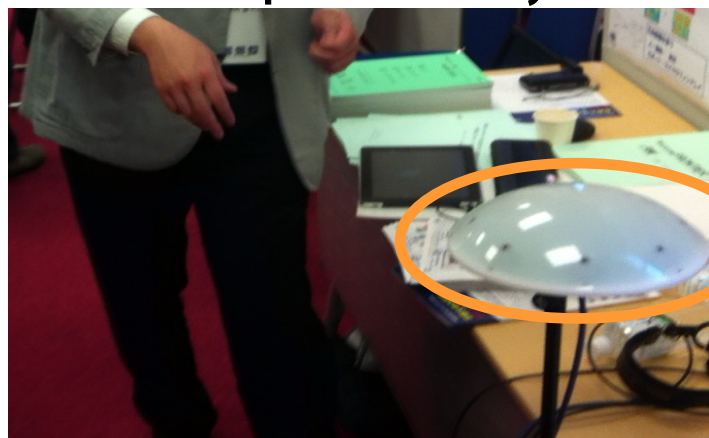
12

- Separate overlapping utterances in a noisy environment

Clock-tower international hall



Microphone array



Observed mixture signal



Separated signals



Background noise



Utterances

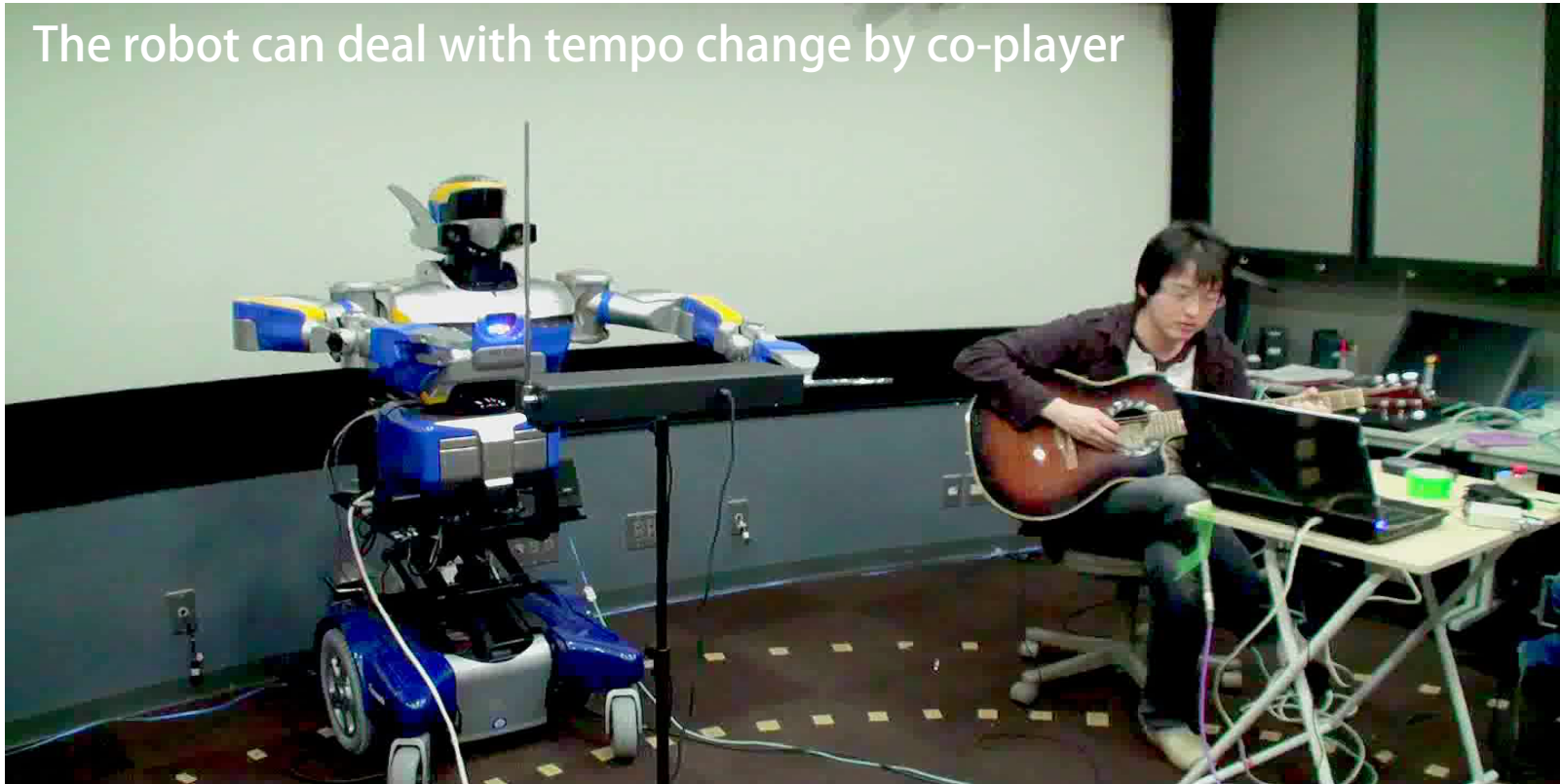
Listen to Music

Music Co-playing with Humans

14

- Real-time score-position tracking
 - Listen to partner's playing by using own ears

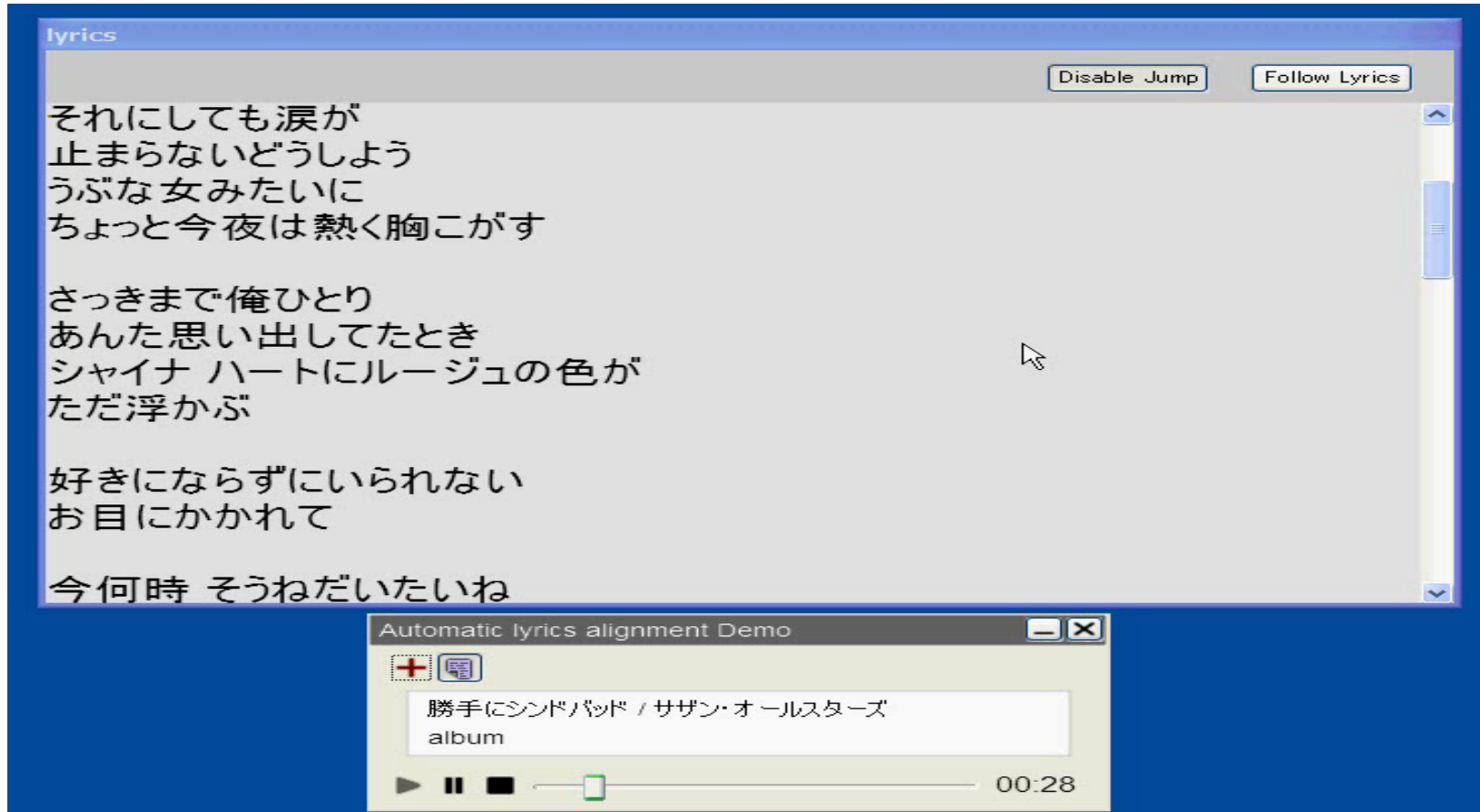
The robot can deal with tempo change by co-player



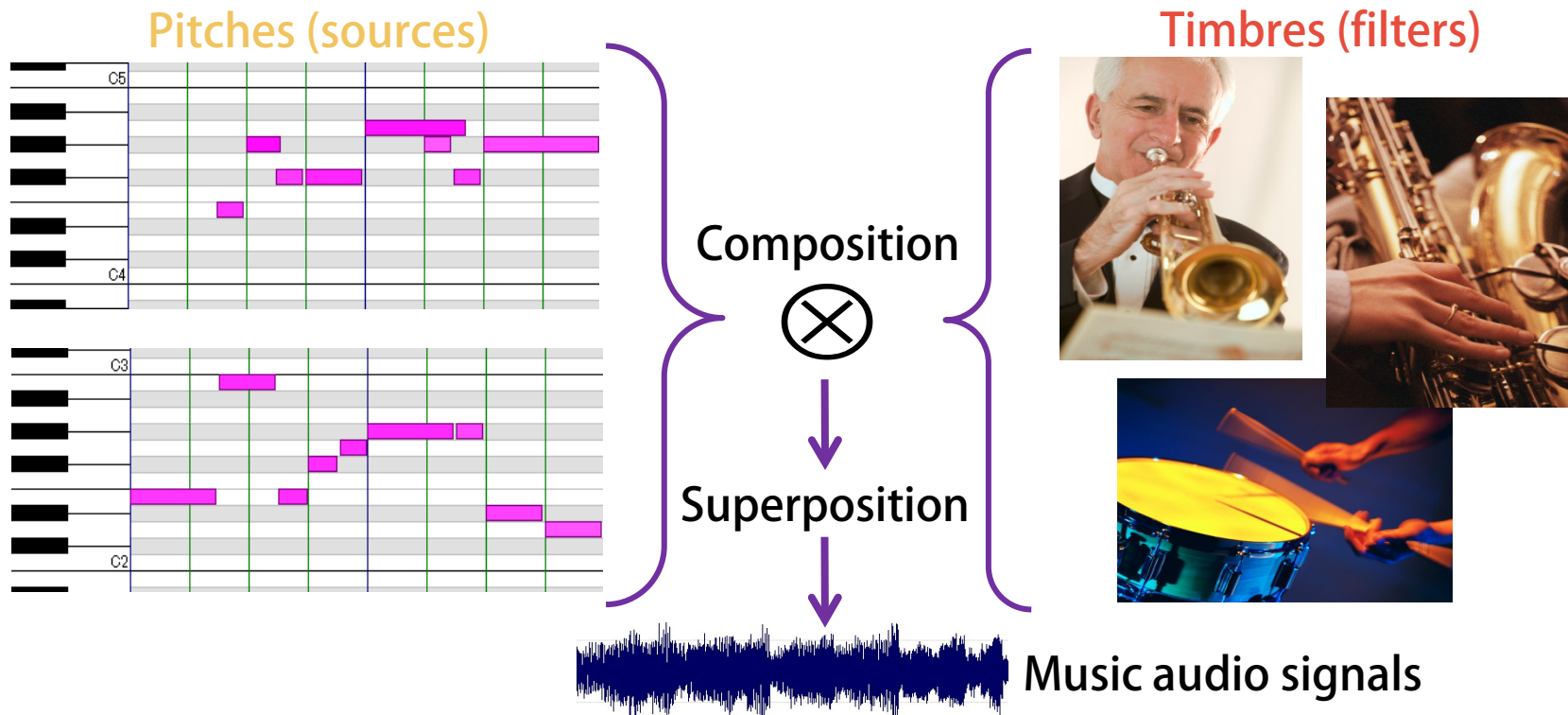
Lyric-to-Audio Synchronization

15

- Efficient navigation to a section of interest



- Parts-based representation of music
 - Combinations of “**pitches**” and “**timbres**”

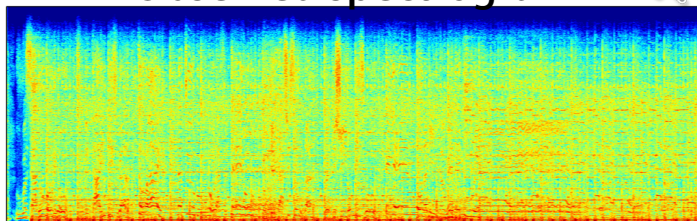


Music Signal Decomposition

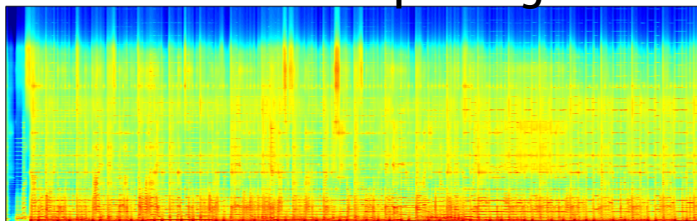
17

- Timbre-based audio source separation

Observed spectrogram



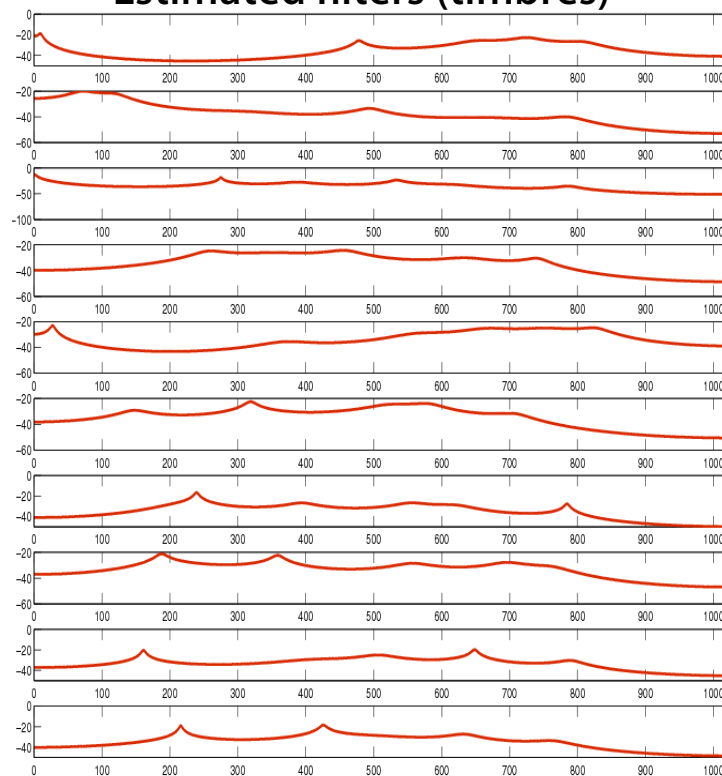
Reconstructed spectrogram



Timbre weights



Estimated filters (timbres)



- Edit only drum parts in mixture signals

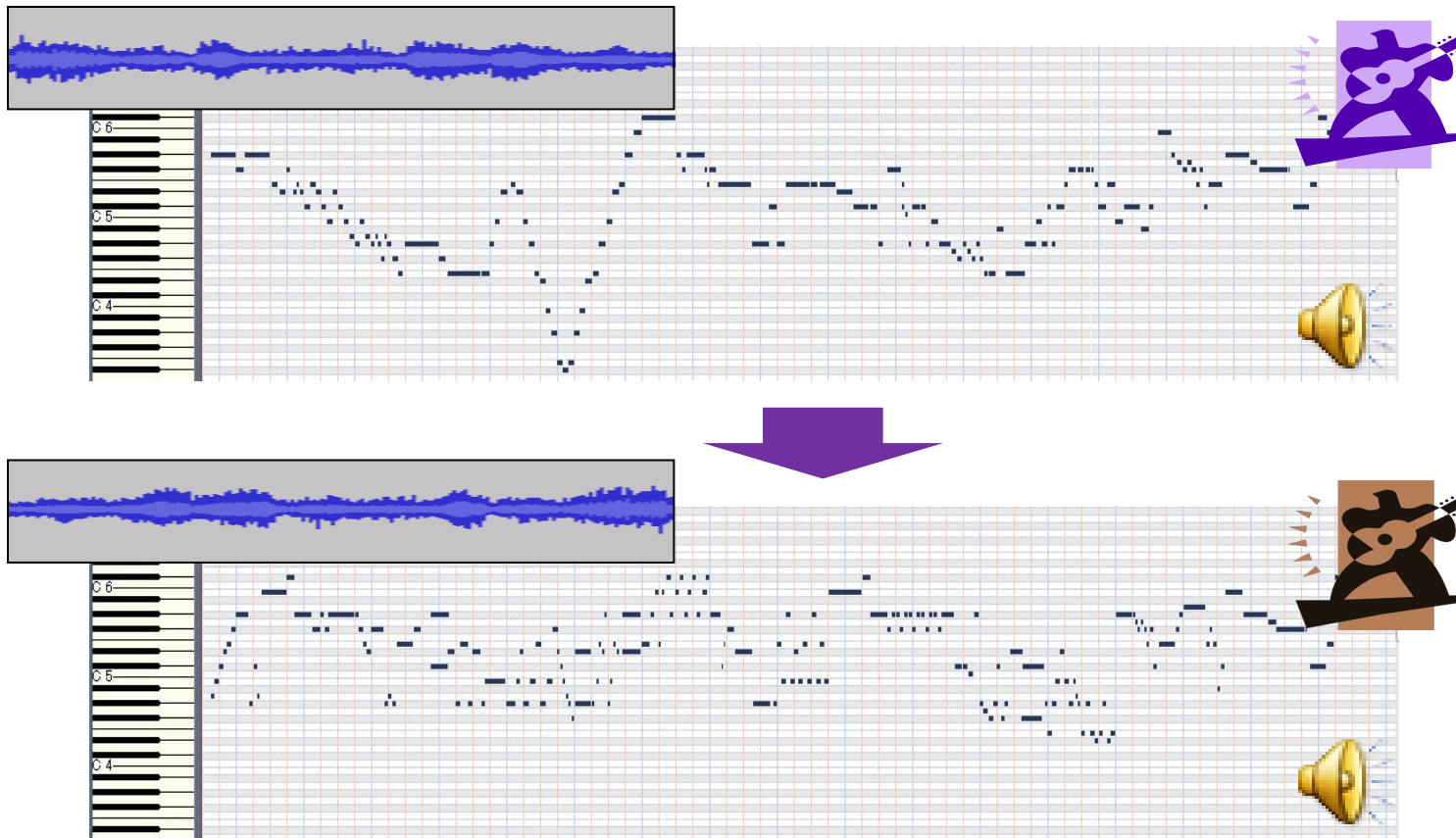
Drumix: An audio player with
a function for re-arranging
drum parts in real time

Kazuyoshi Yoshii
Masataka Goto
Kazunori Komatani
Tetsuya Ogata
Hiroshi G. Okuno

Replacement of Guitar Solo

19

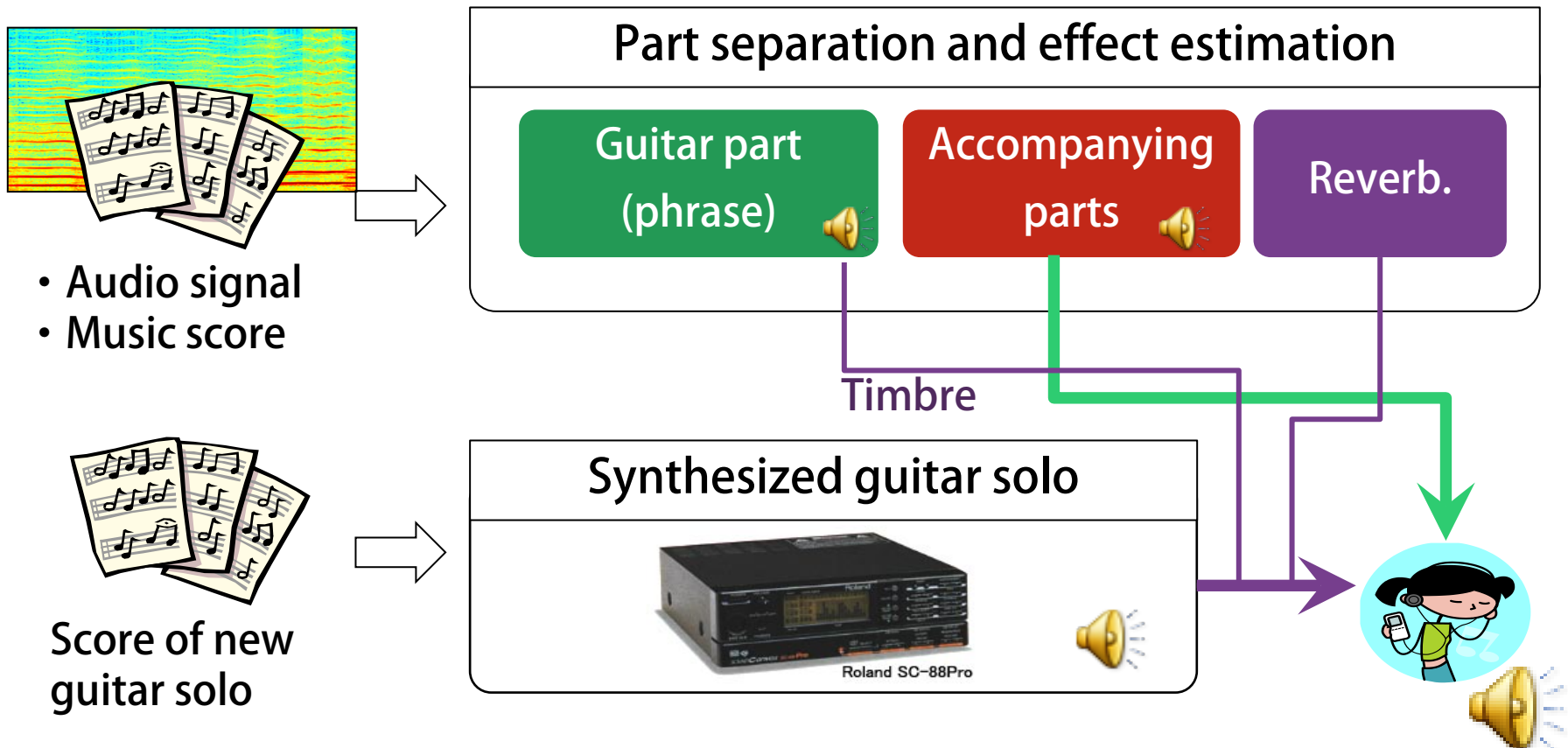
- Edit only a guitar part while preserving original timbres



Timbre and Effect Estimation

20

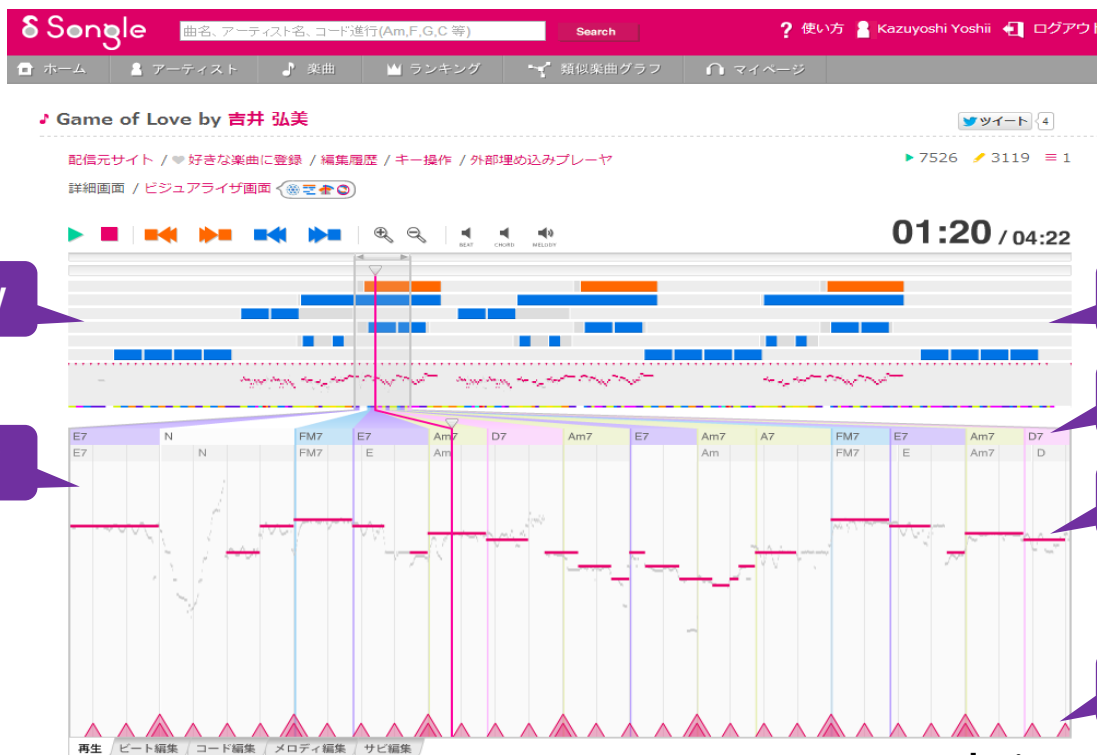
- Preserve timbers and reverberation of original guitar solo



Songle: Active Music listening Service

21

- We can enjoy automatically estimated, visualized, and sonificated musical elements of songs on the Web

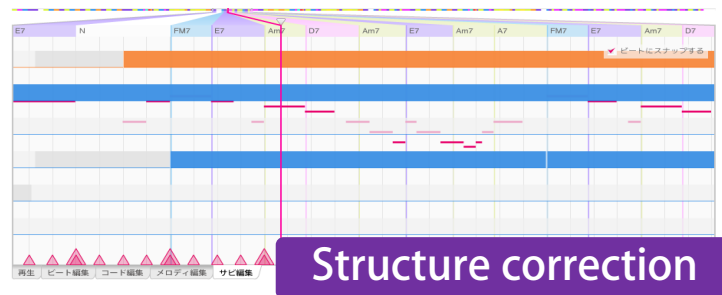
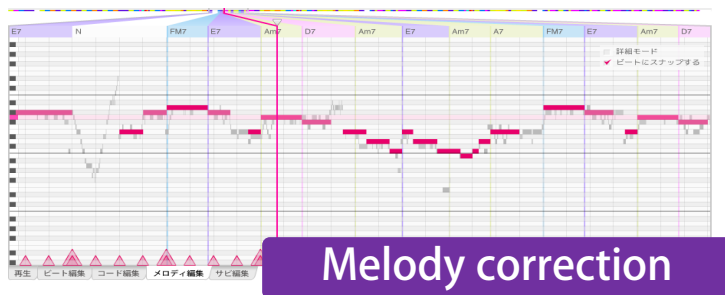
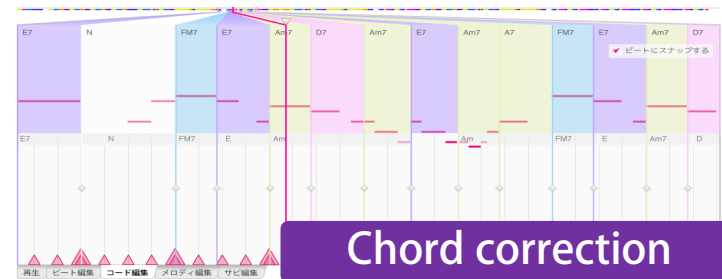
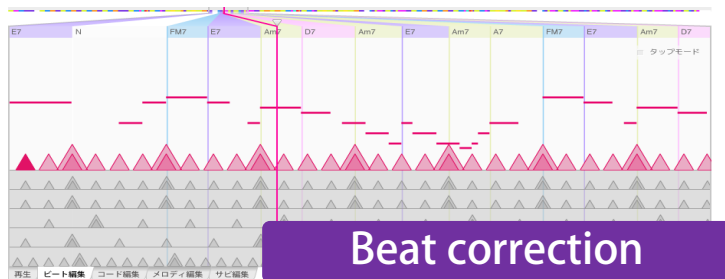


songle.jp

A New Form of Outreach Activity

22

- We can amplify user contributions by using machine-learning techniques
 - Corrections by some users → Retraining → Accuracy improvement
→ **Reward to all users**



Listen to Environmental Sounds

-



Audition in Flying Robots

25

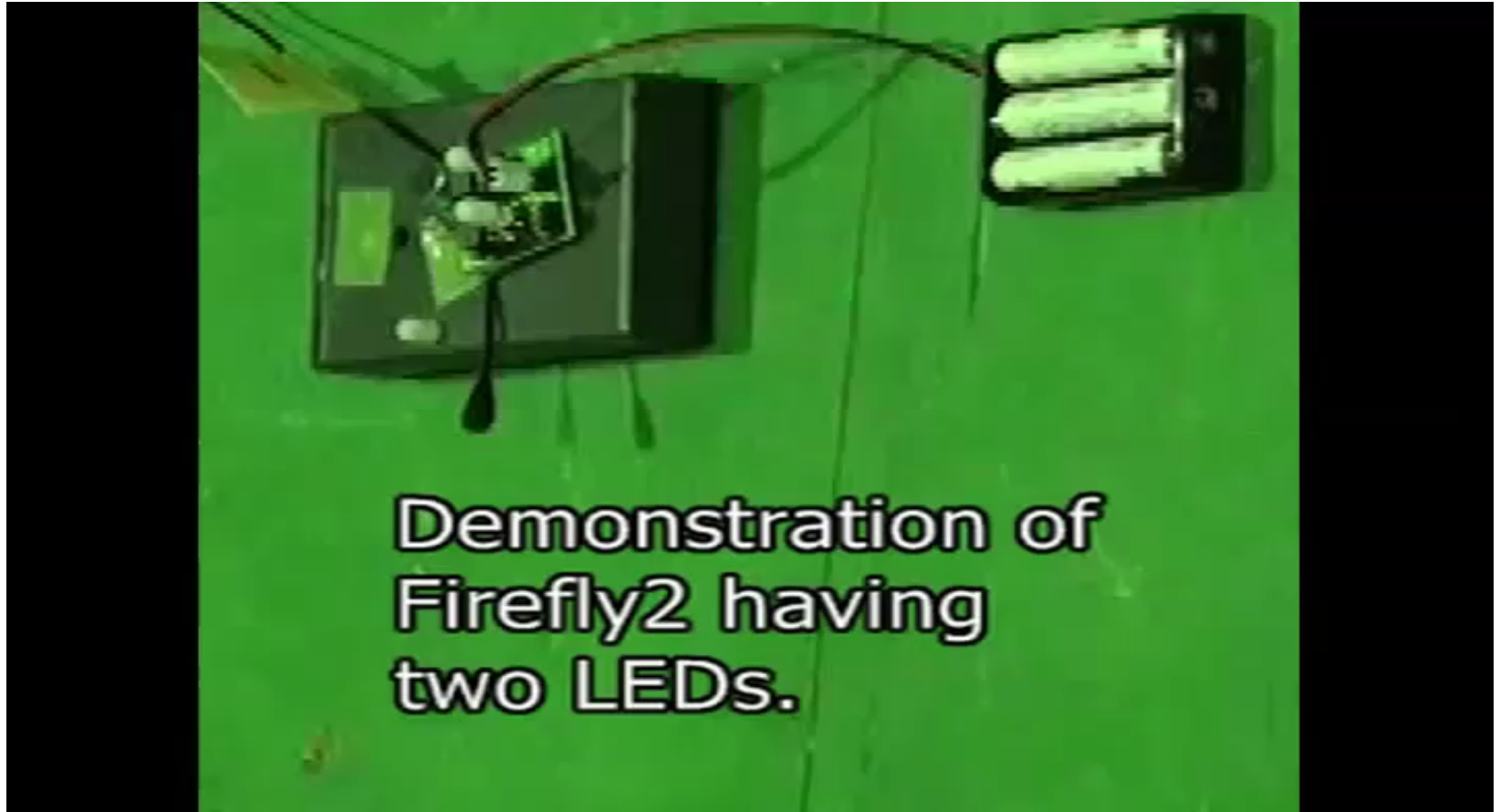
- Use a microphone array for localization



Visualization of Frog Calling

26

- Discriminate calling of two kinds of frogs



- Separation and localization in a park

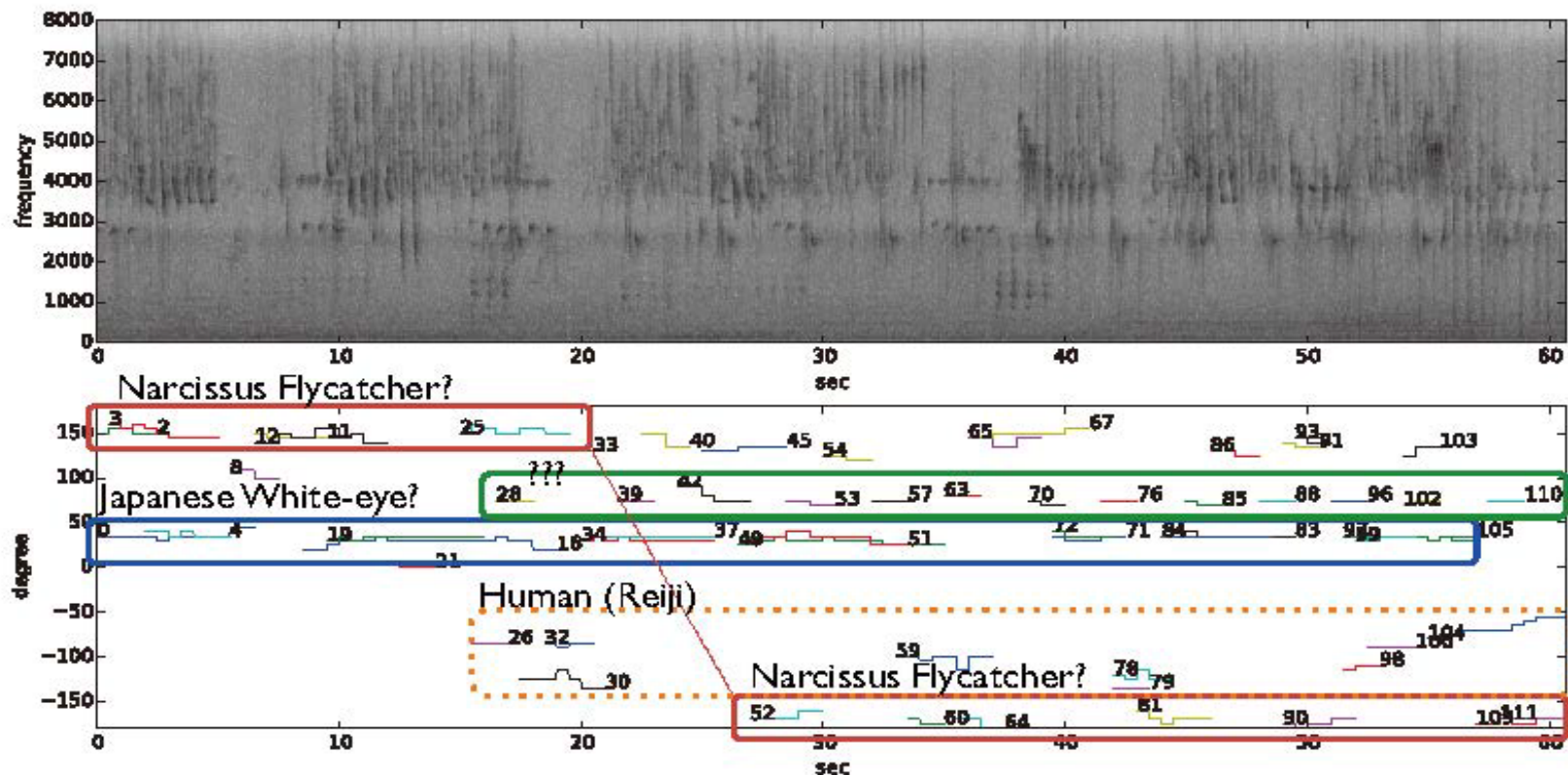


Figure 5: HARK を用いた野鳥の歌の音源定位の例.

Robot Audition

Why Robot Audition?

29



Conventional problem:
We need to speak around microphones

Why?

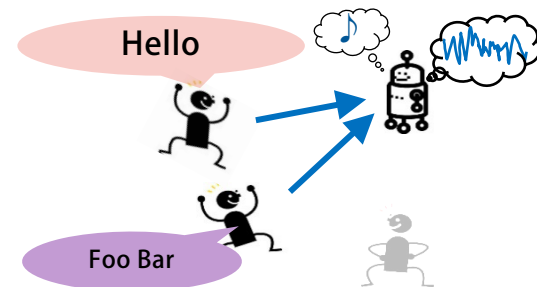
The microphones inevitably catch noise sound with target utterances

Our approach:
We aim to separate and recognize sounds

Sound Source
Localization (SSL)

Sound Source
Separation (SSS)

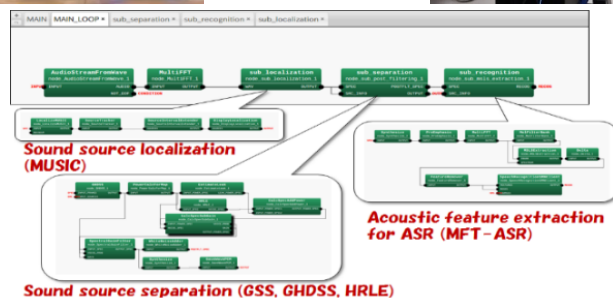
Computational Auditory Scene Analysis (CASA)



- Robot MCs need to interact with multiple people
 - Use an open-source robot audition software “HARK” developed by Honda Research Institute Japan and Kyoto University
 - **HATTACK 25: Speech-based quiz game**

Inspired by the well-known quiz game in Japan

A player position can be identified by his or her voice



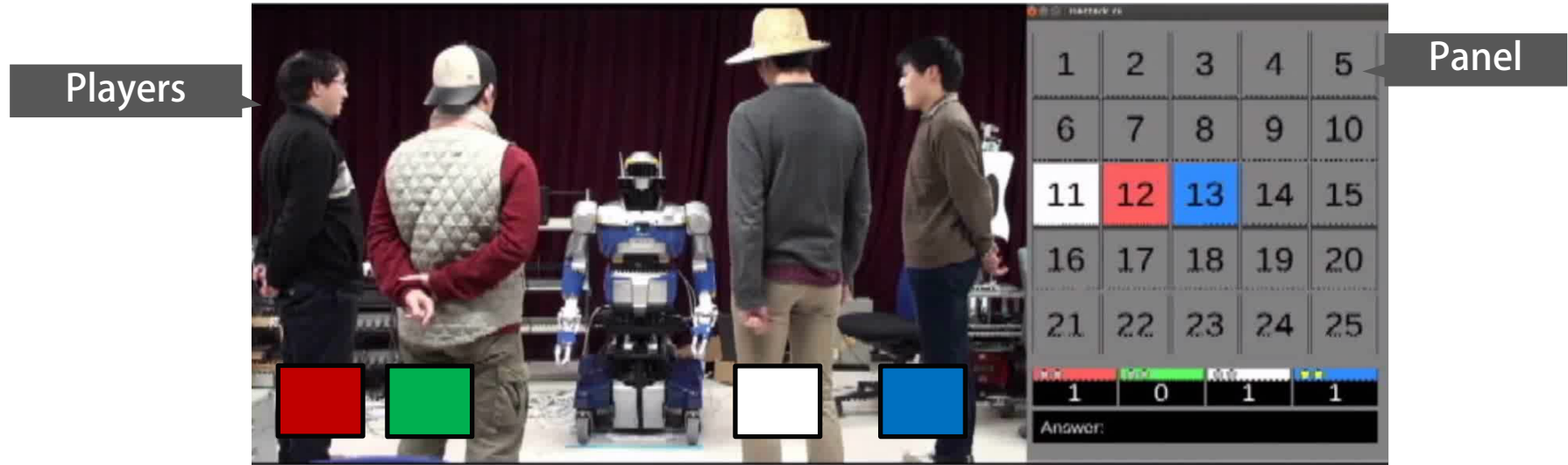
HARK



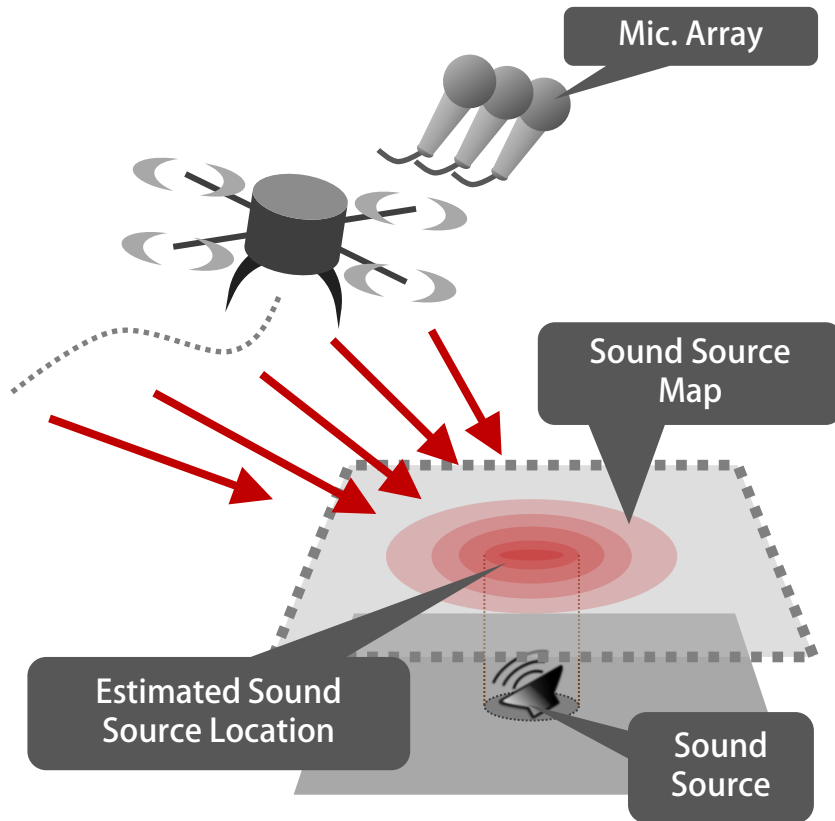
Demo : HATTACK25

31

- Players can barge in when the robot is talking about questions
 - To answer, players have to say “yes!” first
 - Impossible for standard dialogue systems



- Localize source sources on the ground by using flying robots with microphones



Robot audition is disturbed by self-generating noise

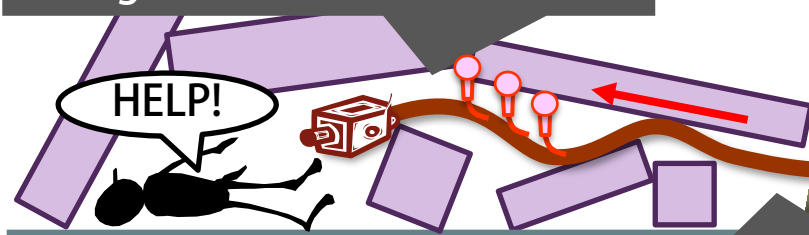
Video of the flying robot



1. Learn self-generating noise (with Gaussian Process)
2. Suppress noise from input

- Estimate robot shapes and detect sound sources in collapsed buildings

Length: 3-8m & Width: 3-5cm



Move forward with self locomotion



Moving hose-shaped robot



How to estimate microphone positions on the robot?
→ Statistical signal processing techniques

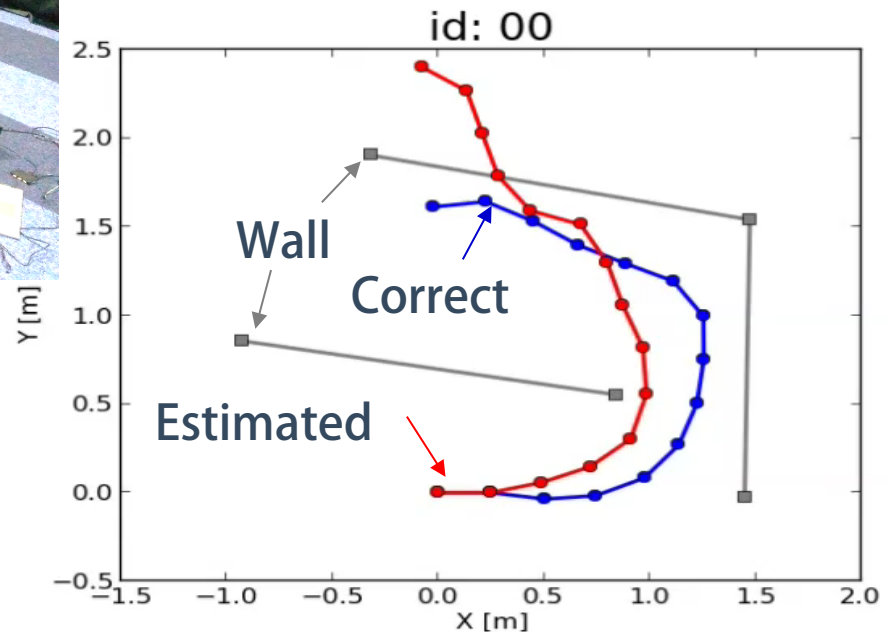
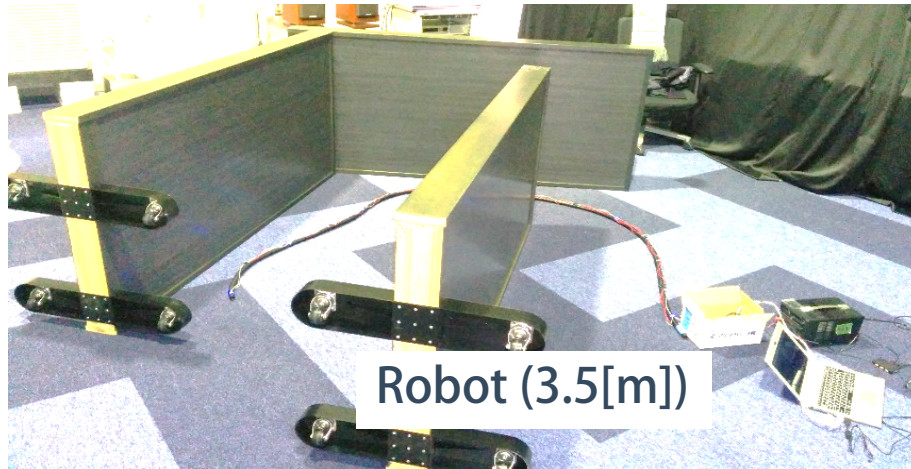
We designed a state space model of robot posture
and estimated the posture by measuring TDOA of sounds

Time Difference of Arrival

Posture Estimation for a Hose-shaped Robot

34

- Accurately estimate shapes even when obstacles exist



- Listen to “speech”
 - Prince-Shotoku robot
 - ♦ Simultaneous speech recognition
 - Microphone-array processing
 - ♦ Sound source separation and dereverberation
- Listen to “music”
 - Music understanding and performance
 - ♦ Sound source separation and music transcription
 - ♦ Co-playing and accompaniment
- Listen to “environmental sounds”
 - Object detection in a disaster environment
 - Analysis of frog calling