

# Robot Audition and Drone Audition

**Kazuhiro Nakadai**

Principal Scientist, Honda Research Institute Japan Co. Ltd.  
Specially-appointed Professor, Tokyo Institute of Technology

**Hiroshi G. Okuno**

Professor, Waseda University  
Professor Emeritus, Kyoto University  
Honorary Professor, Amity School of Engineering and Technology, India

# Robot Audition [Nakadai & Okuno AAI 2000]

HONDA

Honda Research Institute JP

## ■ Not a headset microphone, but *its own ears!*

### – Noise-robustness

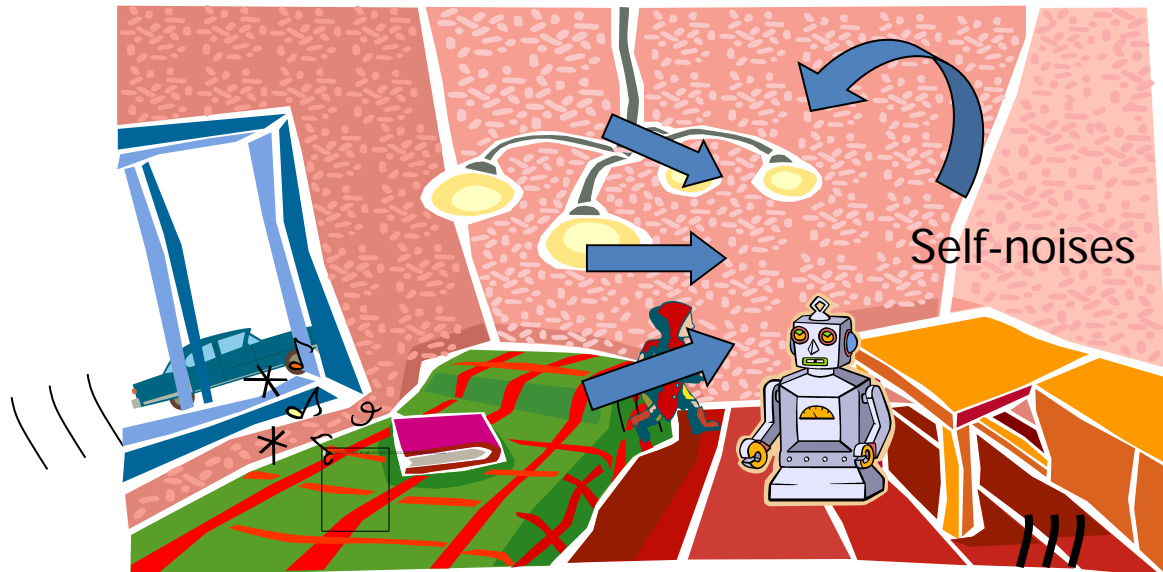
- Ego-noise (actuators, self-voice)
- Environmental sounds
- Simultaneous speech (barge-in)

### – Cocktail Party Robot

### – Prince Shotoku Robot



## ■ Towards Auditory Scene Analysis



# Primary Issues in Robot Audition

## ■ Sound Source Localization (SSL)

- MUSIC based on Generalized Eigen/Singular-Value Decomposition (GEVD/GSVD-MUSIC) [Nakamura+, Okutani+, Ohata+ '09-'12]

## ■ Sound Source Separation (SSS)

- Geometric High-order Decorrelation based Source Separation with Adaptive Step-size Control (GHDSS-AS) [Nakajima+ '10, Takeda+ '12]

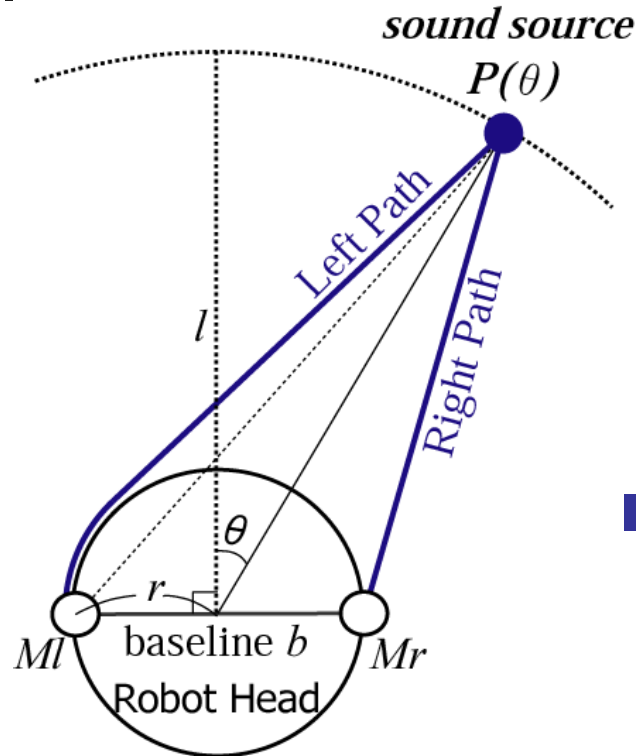
## ■ Automatic Speech Recognition (ASR)

- Missing feature theory based integration of separation and ASR [Yamamoto+ '07]

- Published in multiple research fields; Robotics (ICRA, IROS), Acoustics & Speech (ICASSP, INTERSPEECH), AI (AAAI, IJCAI)
- Organized special sessions, workshop, tutorial, etc
- “Robot Audition” : an official keyword of RAS in 2014



# Auditory Epipolar Geometry



- Head Related Transfer Function (HRTF)
  - Measured in anechoic room
  - prone to alter due to environmental change



- Auditory Epipolar Geometry
  - Method for horizontal localization
  - Estimate direction of a sound source from IPD **mathematically**

$$\Delta\varphi = \frac{2\pi f}{v} \times r (\theta + \sin \theta)$$

(  $\Delta\varphi$  : interaural phase difference (IPD) )

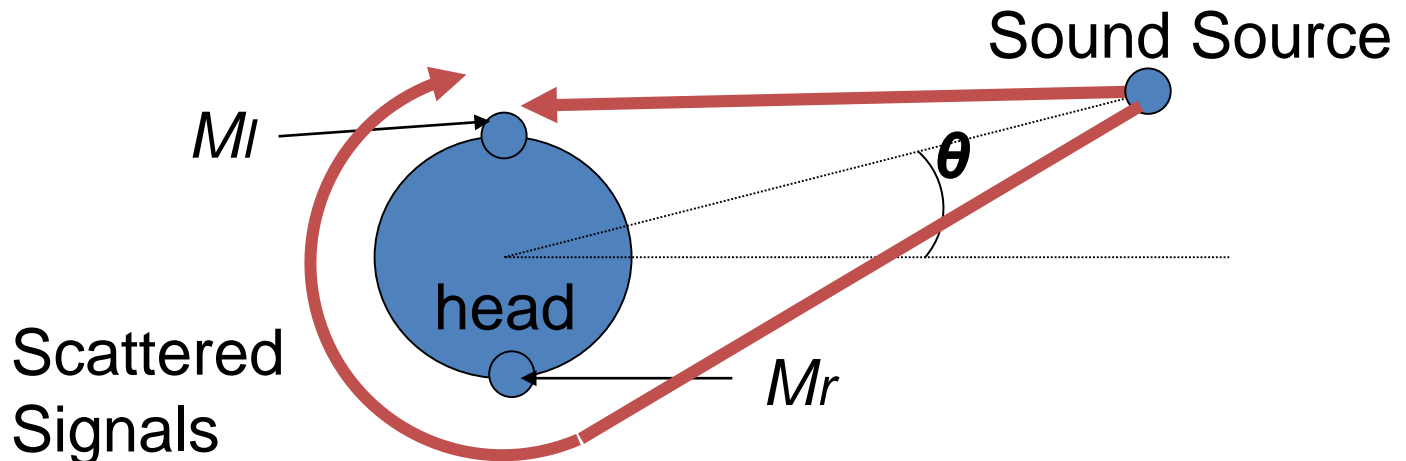
- Based on **Scattering Theory**
- When spherical robot head is assumed, potential at a point on the surface is estimated by

$$S(\theta, f) = - \left( \frac{v}{2\pi a f} \right)^2 \sum_{n=0}^{\infty} (2n+1) P_n(\cos \theta) \frac{h_n^{(1)} \left( \frac{2\pi r_0}{v} f \right)}{h_n^{(1)'} \left( \frac{2\pi a}{v} f \right)}$$

- Estimate IPD and IID **mathematically**

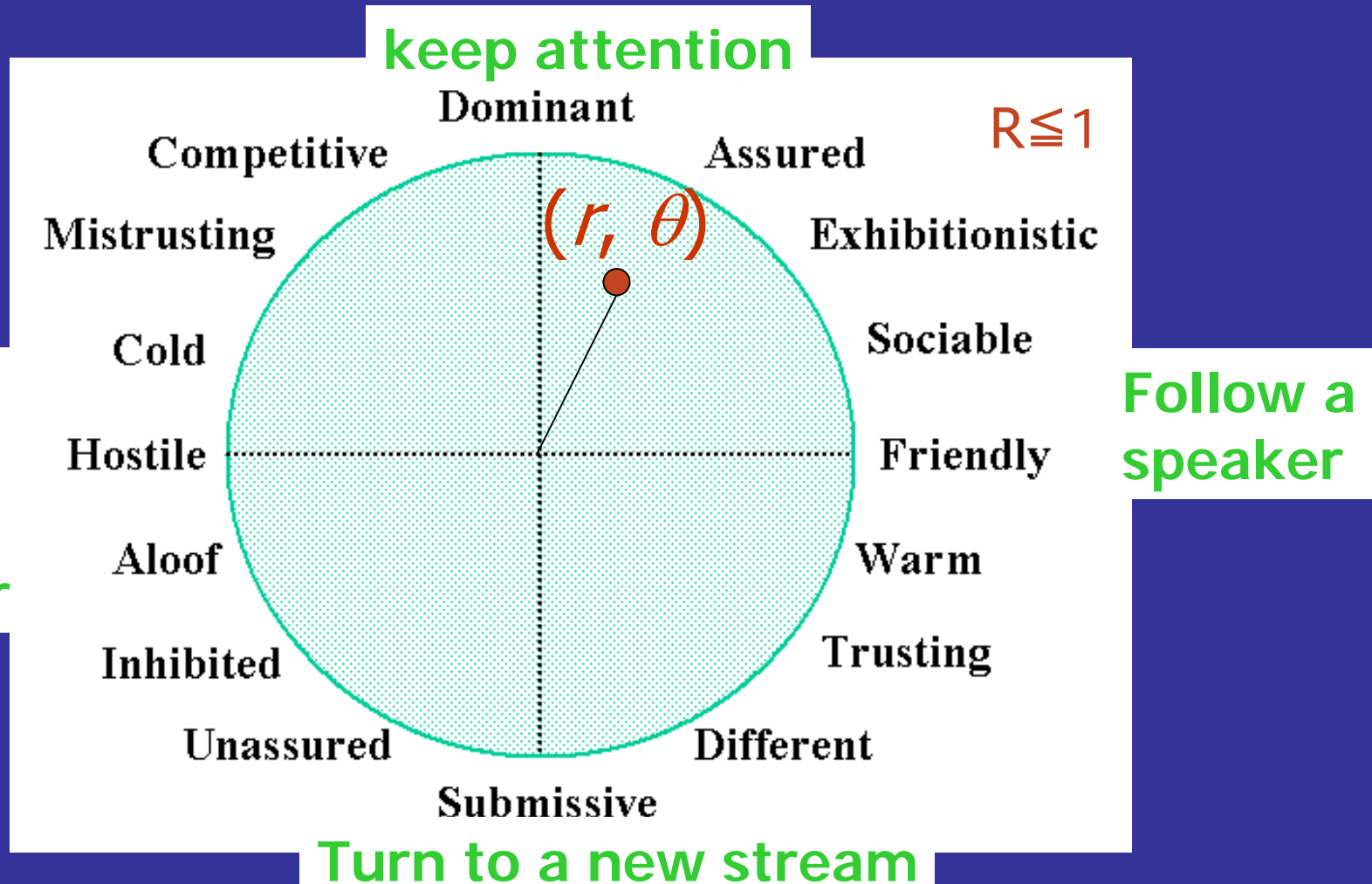
IPD:  $\Delta\varphi_s(\theta, f) = \arg(S_l(\theta, f)) - \arg(S_r(\theta, f))$

IID:  $\Delta\rho_s(\theta, f) = 20 \log_{10} \frac{|S_l(\theta, f)|}{|S_r(\theta, f)|}$



# Personality in Focus-of-Attention [PRICAI 02]

- Personality is represented as a point  $(r, \theta)$  in the Interpersonal Circumplex.



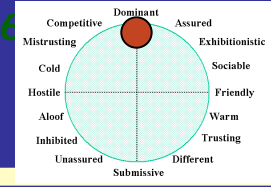
# Demonstrations

---

Demos	Personality ( $r, \theta$ )
1. Follower	Dominant ( $1, \pi/2$ )
<del>2. Receptionist</del>	<del>Dominant (<math>1, \pi/2</math>)</del>
<del>3. Stereo sound tracking</del>	<del>Dominant (<math>1, \pi/2</math>)</del>
<del>4. Companion</del>	<del>Friendly (<math>1, 0</math>)</del>
5. Listening to two people	Assured ( $1, 3/8\pi$ )
6. Hostile listener	Hostile ( $1, \pi$ )
<del>7. lose-interest listener</del>	<del>Friendly with large <math>k</math></del>



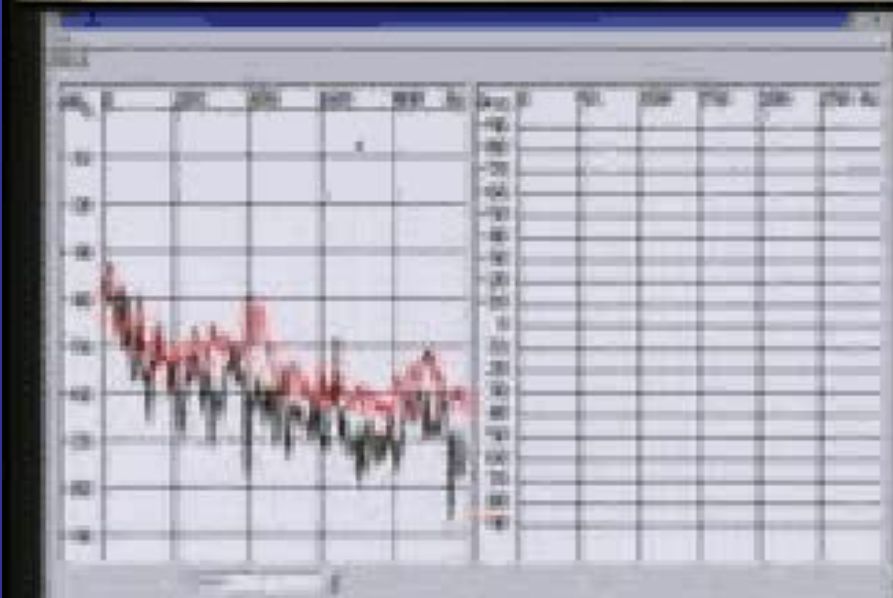
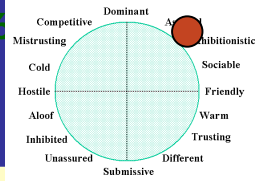
# Follower (Dominant)

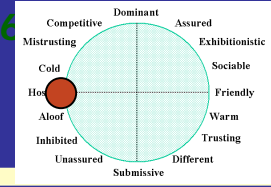




# Listening to two people (Assured)

ICRA-03 (16)





# Hostile Listener (Hostile)



# SIG as a physical non-verbal Eliza

---

1. A variety of non-verbal behaviors is attained by personality-based focus-of-attention control only with sound localization
2. The robot encouraged users to explore the robot's behaviors
  - Some people walk around talking with their hand covering *SIG*'s eyes in order to confirm the performance of auditory tracking.
  - Some people creep on the floor with talking in order to confirm the performance of auditory tracking.
  - Some people play hide-and-seek games with *SIG*.



- **School-type interaction:** ask a right to answer
- **Auction-type interaction:** answer a quiz

## **A Robot Quizmaster for the 'Fastest Voice First' Quiz Game**

Izaya Nishimuta Naoki Hirayama  
Katsutoshi Itoyama Kazuyoshi Yoshii  
Hiroshi G. Okuno

## Simultaneous Speech Recognition

~ Meal Order Taking ~

- Dealing with 11 directional sound sources, a diffuse noise source and ego-noise
- 16ch circular microphone array (speaker locations given).



## ■ HRI-JP Audition for Robots with Kyoto University

(downloadable at <http://www.hark.jp/>)



hark = listen (Old English)

Research: Free  
(Commercial: Licensing)

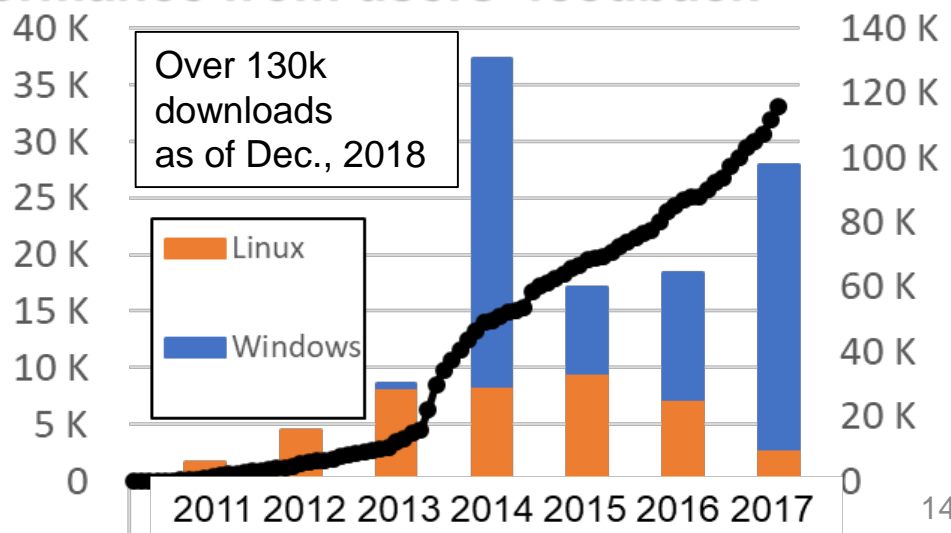
– Provides several algorithms for SSL, SSS & ASR with many utilities.

## ■ Open sourced since Apr. 2008,

- To accelerate robot audition related research
- To provide a tool for collaboration (inter-field, intra-field)
- To improve stability and performance from users' feedback

## ■ Annually update with free tutorial and hackathon

– 15<sup>th</sup> Tutorial was held at IEEE/RSJ IROS 2018





# Research Map of Robot Audition

**ICT/mobility**

Sound search, International communication support, IVI system, UI for optimization, HARK tutorial, Simultaneous Speech recognition, Musical robots, Tele-presence

**Interaction**

HARK cloud service, Sound source separation, Dereverberation, Audio-visual integration, Binaural processing, Interactive dancing

**Robot Audition / Scene Analysis**

Sound Source localization, Noise suppression, Ego noise suppression with non-parametric Bayes model, Sound SLAM

**Ethology/Ecology**

Annotation tools, Outdoor Sound recording, Bird scene analysis (w/ Nagoya Univ.), Onomatopoeia recognition

**Rescue&Search (extreme audition)**

Cyber enhanced canine (w/ Tohoku Univ.), Hose type robots (w/ Kyoto Univ.), UAV sound detection & classification

**Other Projects:** Hearing-aid, Cooking support, Deep Learning, Ring! Ring! Vow Vow

- From Robot Audition & Computational Auditory Scene Analysis to Many engineering & scientific fields
- Base technology: robotics, signal/speech processing, AI incl. DL,...



# What is important in a disaster?

HONDA

Honda Research Institute JP



Great East Japan  
Earthquake '11/03



Kumamoto, Japan  
'16/04



Mexico City, Mexico  
'17/09

- Many disasters (e.g. earthquakes) happen in the world
  - Transportation paralyzed
  - Stacked emergency vehicles
- The Faster, The Better
  - Need to find people within three days to save their lives
    - **Golden 72 hours** in Japan or
    - **The rule of threes** in Western countries
  - One day faster claim = 6 month faster repair (Prof. R. Murphy)



# Why microphone-array-embedded UAV?

HONDA

Honda Research Institute JP

- Realize a new rescue robot, i.e., an **Unmanned Aerial Vehicle (UAV) with a microphone array** to search for people in a disastrous situation
  - **Unmanned Aerial Vehicle (UAV)**
    - UAV *moves quickly over a wide area* even when traffic is cut off.
  - **Microphone array**
    - People can be *detected from acoustic information*, even when people are occluded (a camera is not available).



UAV with a microphone array





# Sound Source Localization: Multiple Signal Classification (MUSIC)

## SEVD-MUSIC [Schmidt 1986]

- Standard Eigen Value Decomposition
- MULTIPLE Signal Classification

Observation:  $x(t)$

FFT

$X(\omega)$

Correlation Matrix

$$R = XX^H$$

$R$

Standard Eigen Value Decomposition

$$R = E\Lambda E^{-1}$$

$E = [e_1, \dots, e_M]$

MUSIC spectrum

$$P(\psi) = \frac{|G(\psi)^H G(\psi)|}{\sum |G(\psi)^H e_m|}$$

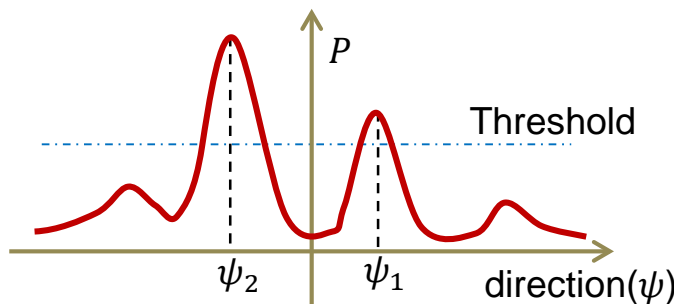
$P(\psi)$

eigenvector  $E$  : sound direction  
 eigenvalue  $\Lambda$  : sound power

Eigenvectors are *orthogonal* from each other.

### Assumptions

- Large eigenvalues : target sound sources  $[e_1, \dots, e_L]$
- Small eigenvalues : noise sources  $[e_{L+1}, \dots, e_M]$



Performance degrades when the target sound source has smaller power than noise sources

MUSIC spectrum

$G(\psi)$  : steering vector (= transfer function) for  $\psi$

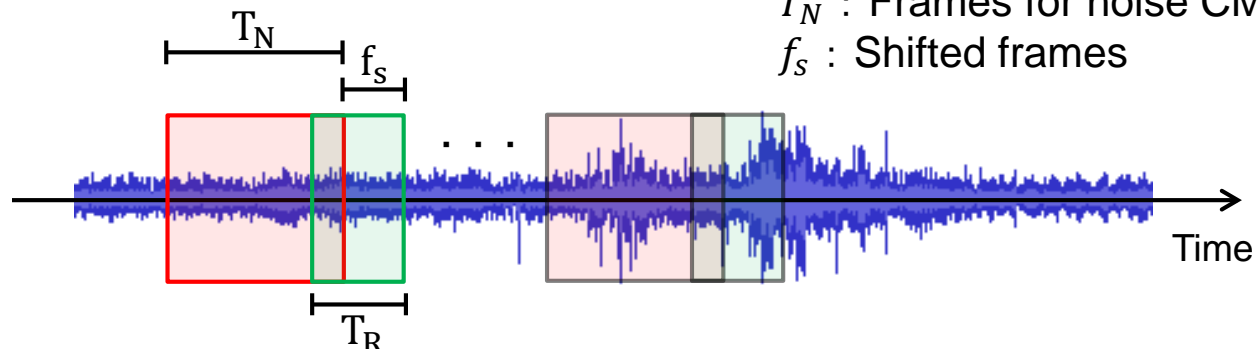
$e_m$  : m-th eigenvectors of  $R$  for noise sources (#of mic. - #of sound sources)

# Extension of MUSIC

## iGSVD-MUSIC [Ohata+'14]

- Whitening noise sources by estimating a noise correlation matrix (CM) incrementally
- To estimate a noise CM, a signal observed  $f_s$  frames before is used by assuming that no target signal is included in the signal.

- for CM,  for noise CM
- $T_R$  : Frames for CM  
 $T_N$  : Frames for noise CM  
 $f_s$  : Shifted frames



Incremental estimation of correlation matrix for iGSVD-MUSIC

A target sound source with smaller power can be localized.

Audio signal

$$x(t)$$

Noise signal

F transpose

F transpose

$$X(\omega)$$

$$N(\omega)$$

Correlation Matrix (CM)

$$R = XX^H$$

Noise correlation matrix (noise CM)

$$K = NN^H$$

$$R$$

$$K$$

Generalized Singular Value Decomposition

$$K^{-1}R = E_L \Lambda E_R$$

$$E_L = [e_1, \dots]$$

Spatial spectrum

$$P(\psi) = \frac{|G(\psi)^H G(\psi)|}{\sum |G(\psi)^H e_m|}$$

$$P(\psi)$$

$G(\psi)$  : Spatial transfer function  
 $e_m$  : Singular vectors for noise sources





# Offline Sound Source Localization with iGSVD-MUSIC

[Ohata+14, Nagamine+14]

HONDA

Honda Research Institute JP

地面レベルのビューを終了

(412.	82.4.	19.338572,	-80'	-115')
(413.	82.6.	19.377977,	-85'	-120')
(414.	82.8.	19.357016,	-85'	-120')
(415.	83.0.	19.124084,	-15'	165')
(416.	83.2.	19.132011,	-60'	-35')
(417.	83.4.	19.108383,	-55'	-35')
(418.	81.6.	19.135944,	-45'	80')
(419.	83.8.	19.134338,	-10'	90')
(420.	84.0.	19.148174,	0'	35')
(421.	84.2.	19.13723,	-45'	-175')
(422.	84.4.	19.159948,	-55'	20')
(423.	84.6.	19.127758,	-65'	-125')
(424.	84.8.	19.172302,	-20'	-90')
(425.	85.0.	19.138767,	-75'	45')
(426.	85.2.	19.150257,	-45'	-135')
(427.	85.4.	19.185745,	-35'	-120')
(428.	85.6.	19.133379,	-35'	-120')
(429.	85.8.	19.098965,	-30'	35')
(430.	86.0.	19.144064,	-85'	60')
(431.	86.2.	19.191576,	-30'	-25')
(432.	86.4.	19.224384,	-75'	-115')
(433.	86.6.	19.145874,	-70'	165')
(434.	86.8.	19.141886,	-35'	-145')
(435.	87.0.	19.128067,	-50'	-25')
(436.	87.2.	19.12566,	-70'	-170')
(437.	87.4.	19.133251,	-65'	180')
(438.	87.6.	19.122759,	-15'	30')
(439.	87.8.	19.120035,	-75'	125')
(440.	88.0.	19.185553,	-65'	-120')
(441.	88.2.	19.208746,	-65'	-120')
(442.	88.4.	19.177656,	-65'	-120')
(443.	88.6.	19.171312,	-75'	-100')
(444.	88.8.	19.130213,	-80'	180')
(445.	89.0.	19.177242,	-50'	-135')
(446.	89.2.	19.169945,	-50'	-140')
(447.	89.4.	19.163668,	-50'	-140')
(448.	89.6.	19.123415,	-35'	50')
(449.	89.8.	19.114365,	-45'	-165')
(450.	90.0.	19.115595,	-45'	-165')
(451.	90.2.	19.123955,	-25'	-25')
(452.	90.4.	19.114321,	-40'	-160')
(453.	90.6.	19.148535,	-65'	-115')
(454.	90.8.	19.198076,	-35'	-120')
(455.	91.0.	19.170452,	-40'	-125')
(456.	91.2.	19.204926,	-35'	-120')
(457.	91.4.	19.16864,	-60'	25')

Image © 2014 DigitalGlobe

Sound captured with Video Camera

Google Earth

ツアーガイド 2012 画像取得日: 2013/4/15 高度: 1.11 km

Extremely noisy sound sources (-15dB) were successfully localized.



# Issues and approaches for online demo

HONDA

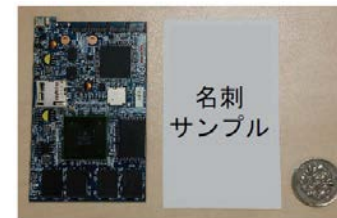
Honda Research Institute JP

## 1. Real-time processing, communication reduction

- Achieved real-time processing and reduction of communication to 1/100 developing an embedded microphone array system

## 2. 3D sound source localization

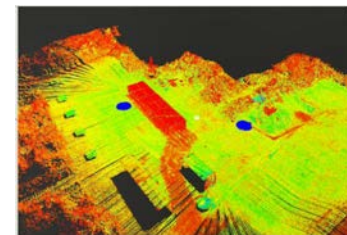
- Estimate 3D position by triangulation with weighted LMS [Washizaki+ IROS 2016]



RASP-MX  
(SiF, Inc)

## 3. Robustness for online outdoor demo

- Take "less time to get ready" with "less human errors" in "all weather"



Real-time visualization  
Blue points are sound locations

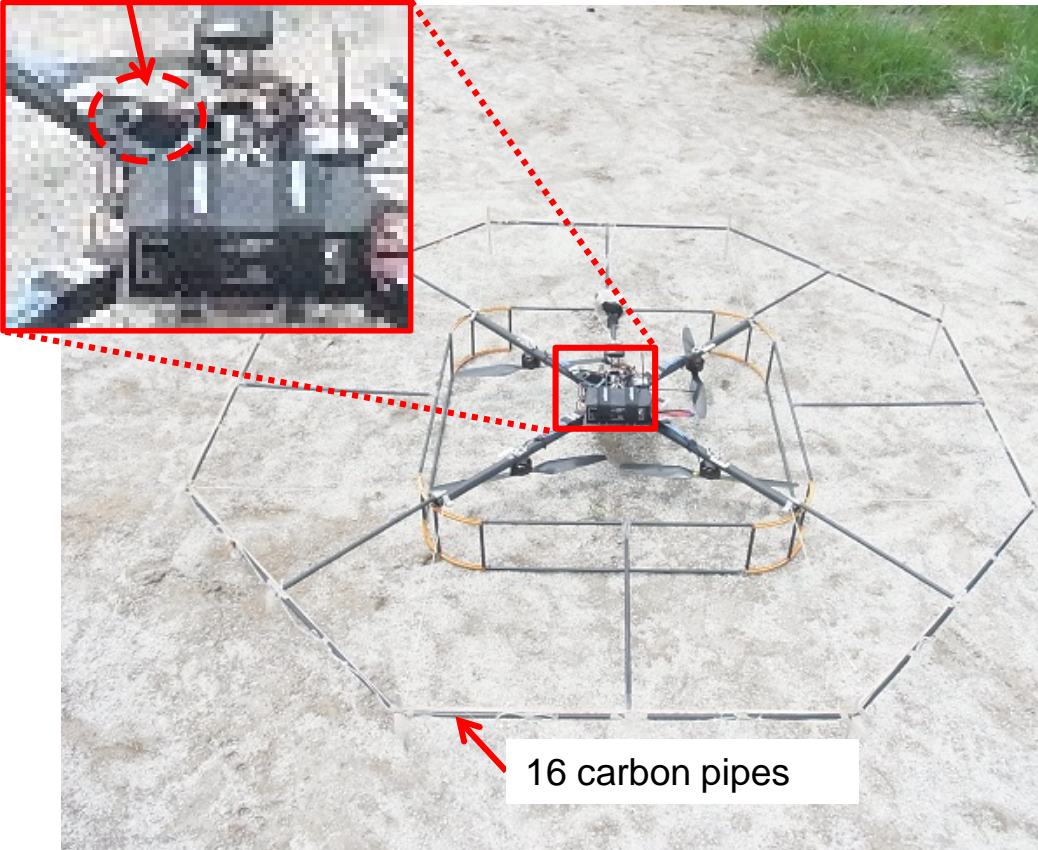
- Intuitive UI
  - 3D visualization on a point cloud map
- Microphone array structure
  - 16ch spherical array
- All weather
  - Water-resistant microphone array



Water-resistant test (passed 12h submersion test)

# First UAV system for online demo

RASP-ZX  
(naked)



## Prototype: enRoute UAV

- 16 pipes and many cables
- Octagonal layout of microphones
- Naked multichannel audio device RASP-ZX for cooling purposes



- *2-hour setup time* necessary
- Many *disconnections* and *contact failures*
- *Not water-resistant*

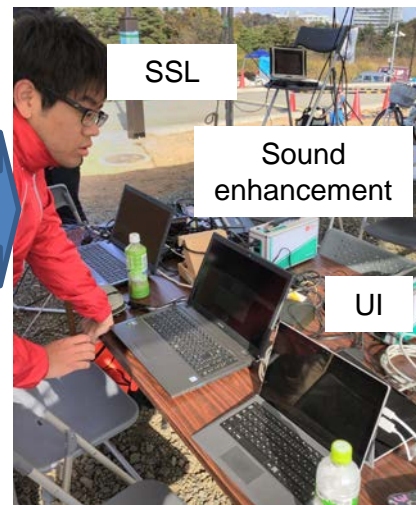
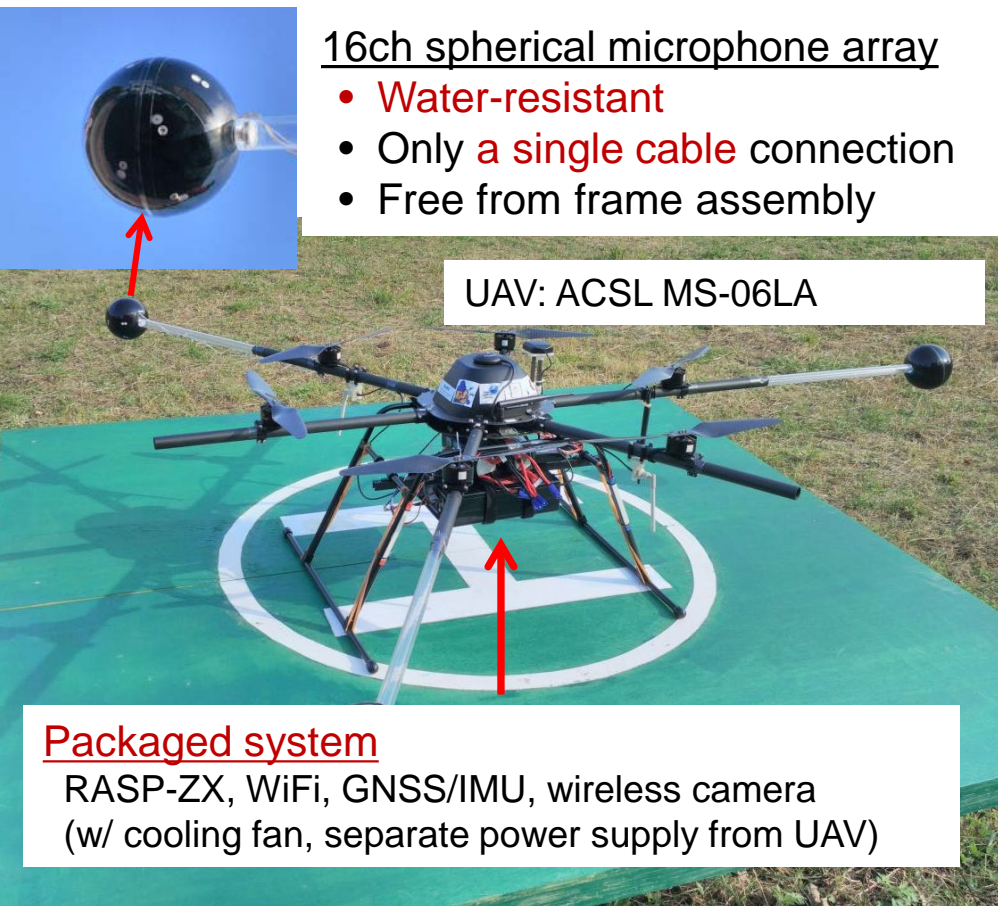


# Revised UAV system for online demo

[Nakadai+ IROS 2017]

HONDA

Honda Research Institute JP



Parallel distributed processing with HARK and ROS



Preparation Time: **reduced 2 hours → 40 min**

Just 2min necessary to take off after the UAV is switched on.

Human errors/disconnections: **not observed**

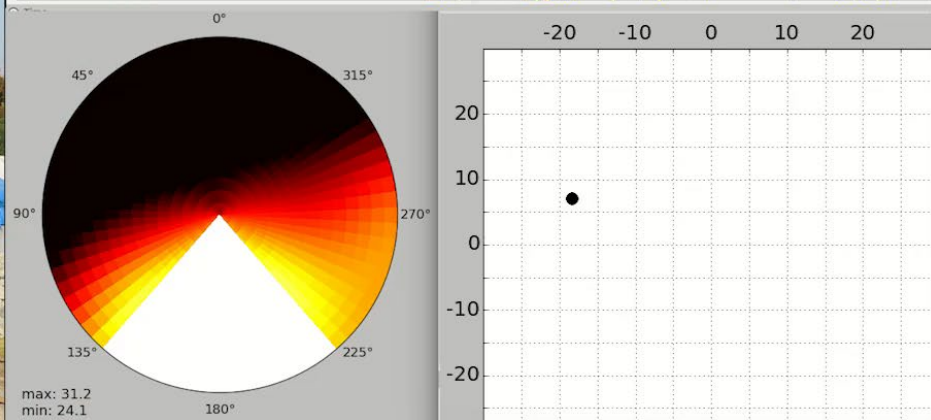
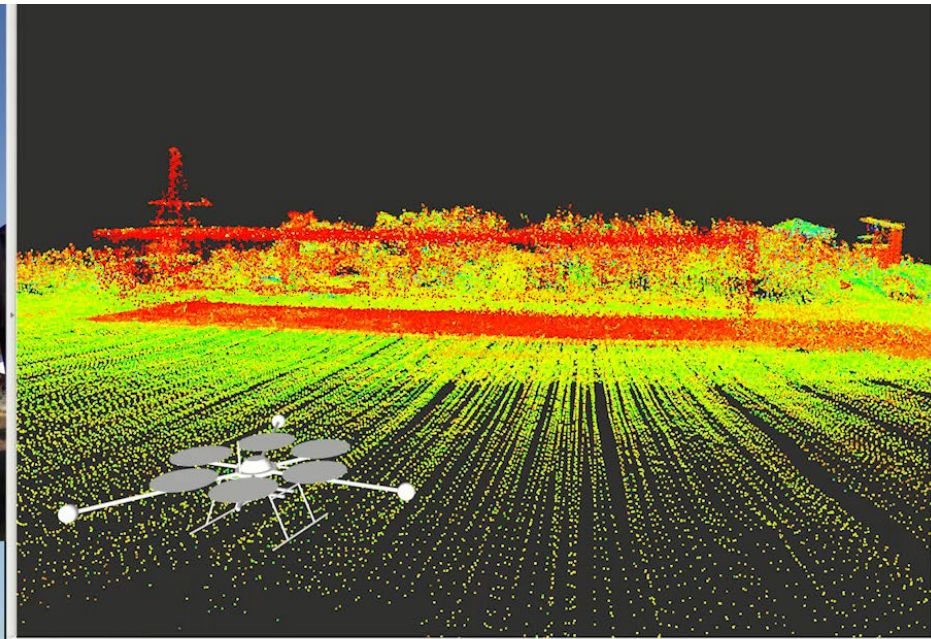
Weather conditions: **operates even in rainy conditions**



# Online Demo@Nov. 11, 2017 (ImPACT TRC field evaluation)

**HONDA**

Honda Research Institute JP

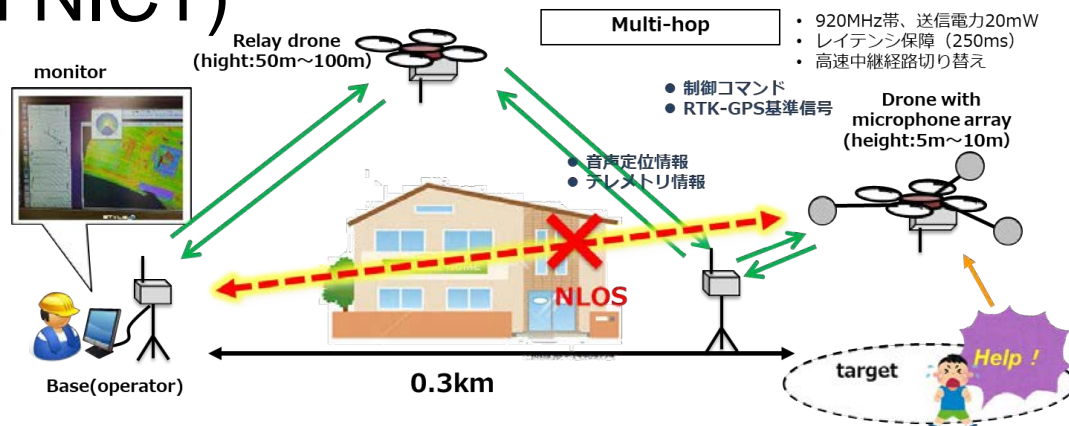


Detected sound locations, and displayed them on 3D map

# For more practical tasks

## ■ Non Line of Sight (NLOS)

- Use multi-hop wireless communication (collaboration with NICT)

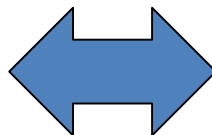


## ■ High Mobility

- Slow speed with arms (microphone array) spread
- An open & close mechanism installed



Motion mode (arm closed)



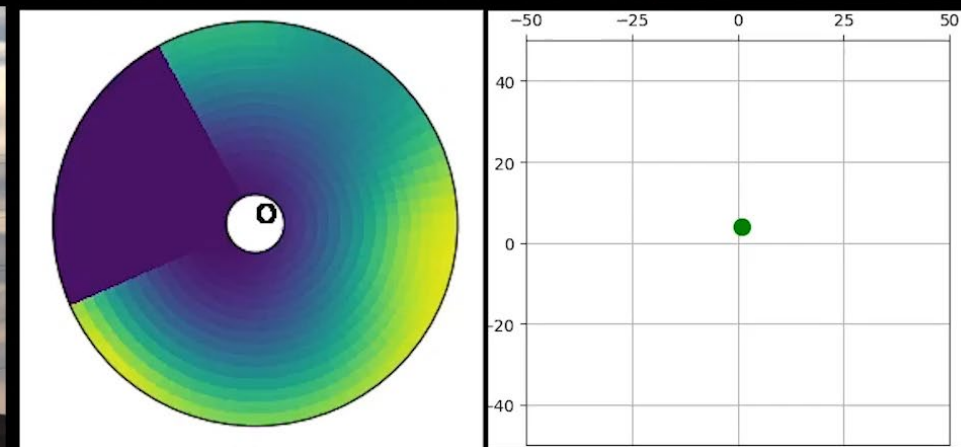
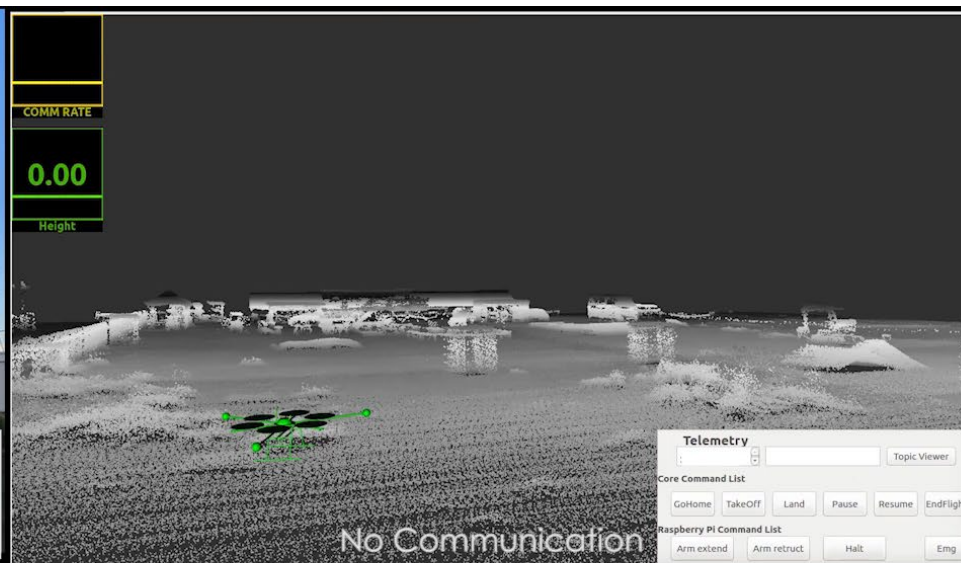
search mode (arm open)



# Final Demo@Nov. 2, 2018 (ImPACT TRC field evaluation)

HONDA

Honda Research Institute JP



Copyright (c) 2018 Tokyo Tech, Kumamoto Univ. Waseda Univ. NICT, AIST (ImPACT TRC demo)

# Sound Source Localization Based on Deep Learning [Nelson+ 2016]

## • Deep Residual Network [He+ 2015]

- Winner algorithm for ILSVRC 2015 based on CNN
- Making optimization easier by learning residue
- Easy to increase the number of layers by avoiding vanishing

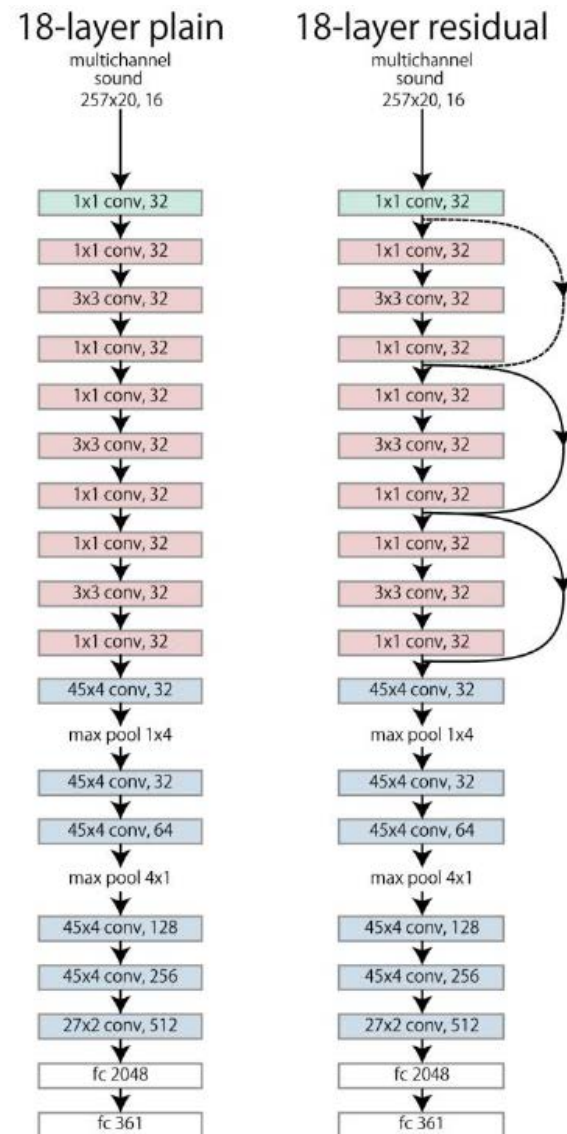
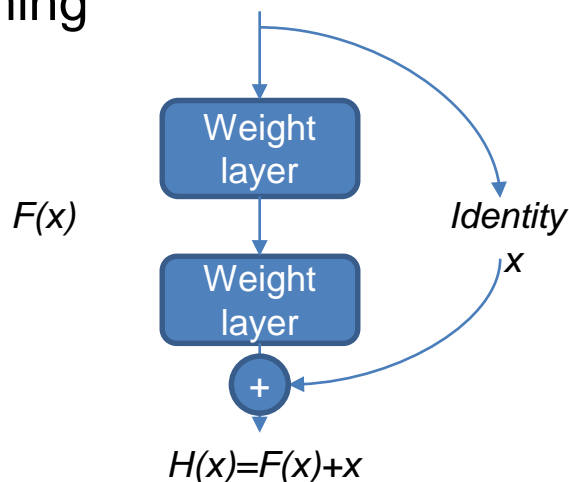
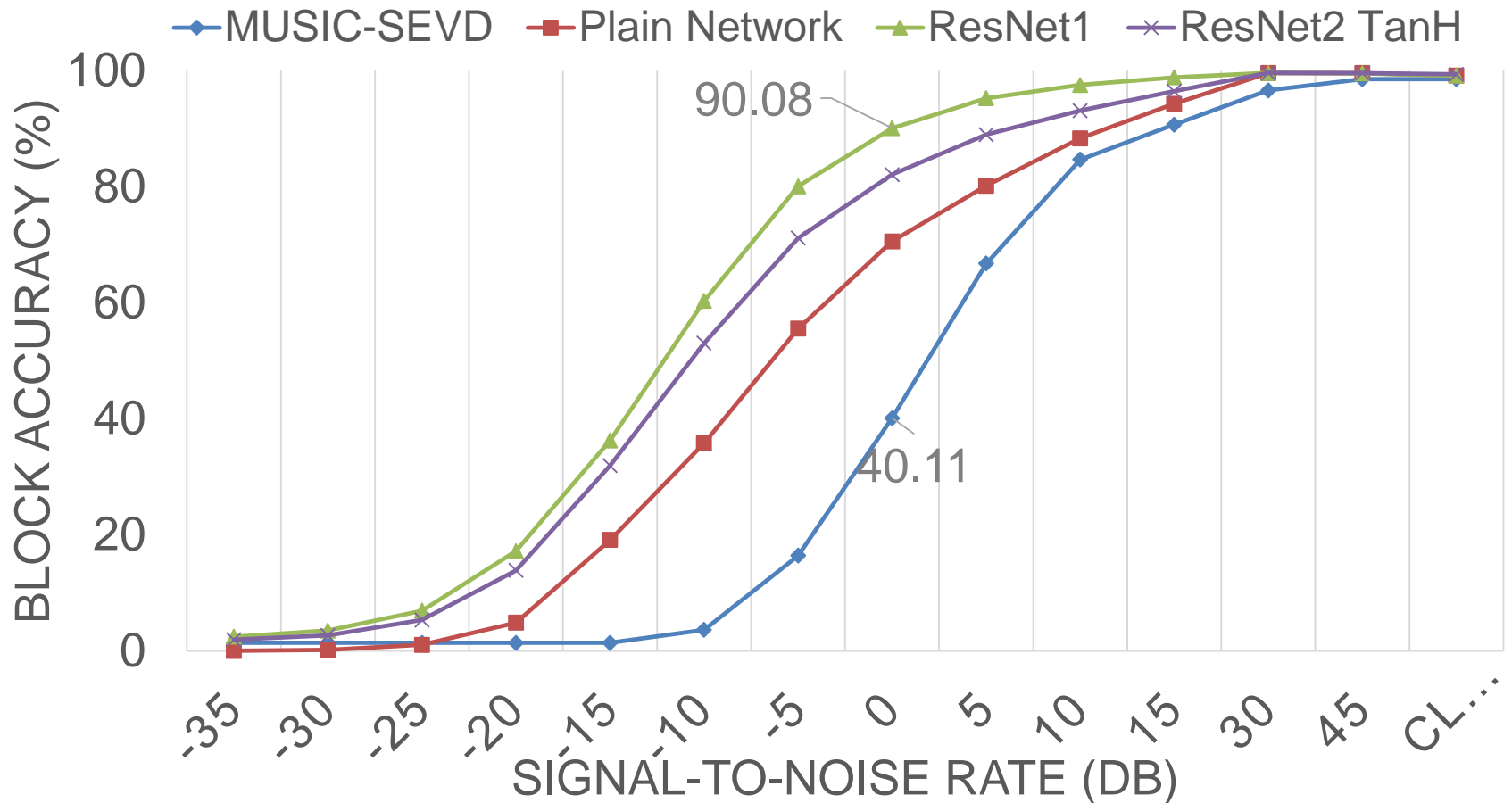


Fig. 4 Network Architecture

Left: Plain Network Right: Residual Network. The dotted line shortcut represents a 1x1 convolutional layer 27

# Performance of Sound Source Localization



- Training data: JNAS (3,350,000 files),
  - Add robot's motor noise to training data(SNR clean ~ -35dB)
  - Input: 257-dim STFT feature (raw)
  - Output: 36-dim localization result (10 deg resolution)
- Test data: 7,200 files ( 200 per direction)

Amplitude input only, and no phase information



# Summary

- **Introduced robot audition and drone audition**
- **Sound source detection for a UAV with a microphone array**
  1. Real-time processing, communication reduction
  2. 3D sound source localization
  3. Robustness for online outdoor demo
  4. NLOS wireless communication
  5. Open&close microphone array for high mobility
- **Our activities in ImpACT (finished Mar, 2019) were published a book Disaster Robotics 2019**
- **Future work**
  1. Frame-based sound source localization
    - Multiple microphone arrays [Yamada+ 2019] workshop poster
  2. Sound source classification [Morito+16, Uemura+17, Nakadai+18]
  3. Multimodal scene understanding: visual processing, thermal camera, etc
  4. Deep sound source localization for dynamic environments