

ロボット聴覚研究とその展開

～災害時の迅速な要救助者発見に向けた
ドローン聴覚技術開発に至るまで～

中臺 一博

(株) ホンダ・リサーチ・インスティテュート・ジャパン プリンシパル・リサーチチャ
東京工業大学 工学院 システム制御系 特任教授

ロボット聴覚 [AAAI 2000]

■ ヘッドセットではなく、自らの耳で！

－ ノイズ頑健性

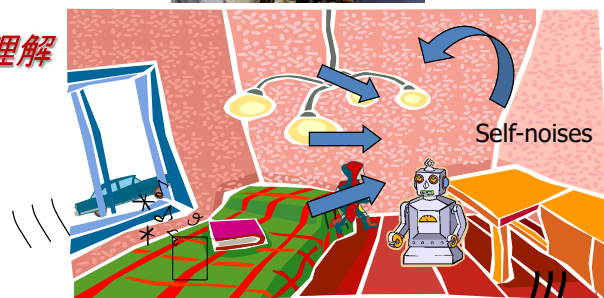
- ・ 自己雑音（モータ、自分の声）
- ・ 環境雑音
- ・ 同時発話（バージン）

－ カクテルパーティロボット

－ 聖徳太子ロボット



■ 何が雑音？⇒環境理解 (Scene Analysis) の 必要性



本日の発表

1. ロボット聴覚の紹介

－ *HARK* の概要・特長

2. ロボット聴覚からドローン聴覚への展開

- － 雑音下音源探索
- － 雑音下音源同定

3. 深層学習へのシフトと課題

4. まとめ

SFにみるロボット (1950)



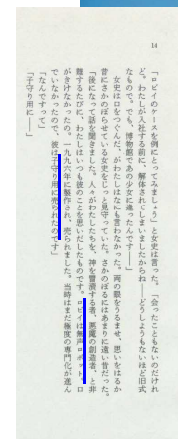
"Take the case of Robbie," she said. "I never knew him. He was dismantled the year before I joined the company—hopelessly out-of-date. But I saw the little girl in the museum—"

She stopped, but I didn't say anything. I let her eyes mist up and her mind travel back. She had lots of time to cover.

"I heard about it later, and when they called us blasphemers and demon-creators, I always thought of him. Robbie was a non-vocal robot. He couldn't speak. He was made and sold in 1996. Those were the days before extreme specialization, so he was sold as a nursemaid—"

"As a what?"

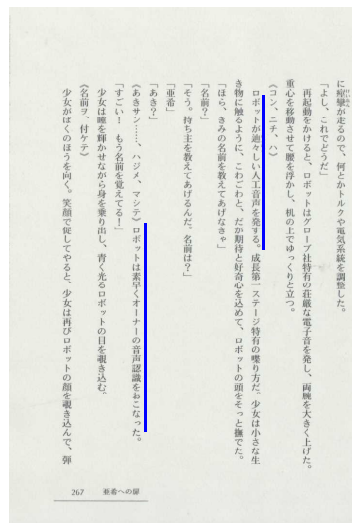
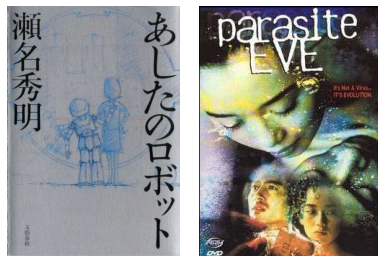
"As a nursemaid—"



最初のロボット "Robbie" (1996 発売)

- ・ Nursemaid → 人の話を理解できる.
- ・ Non-vocal → しゃべれない.

SFにみるロボット (2002)



■ 拾ってきたロボット「Robbie」

- 音声認識: 素早く, 正確
- 音声合成: たどたどしい

6

ロボット聴覚の主要課題

■ 音源定位 (Sound Source Localization)

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

■ 音源分離 (Sound Source Separation)

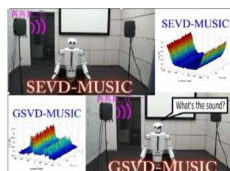
- Geometric High-order Decorrelation based Source Separation with Adaptive Step-size Control (GHDSS-AS) [Nakajima+ '10]

■ 音声認識 (Automatic Speech Recognition)

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

ロボット用のアルゴリズムの研究開発を推進

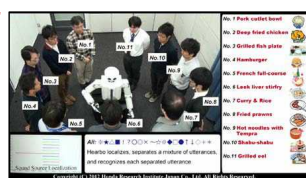
- ロボット (ICRA, IROS), 音響, 音声 (ICASSP, Interspeech) AI (AAAI, IJCAI) の国際会議を中心に発表



雑音ロボスト音源定位



4人同時発話分離



11人同時発話認識

8

ロボット聴覚の研究としての位置づけ

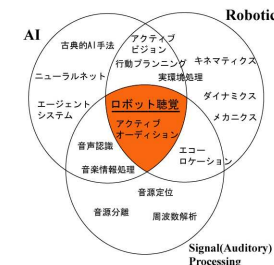
■ 日本発の研究分野: 奥乃・中臺 [AAAI-2000]

- <http://winne.kuis.kyoto-u.ac.jp/SIG/>

■ 工学的には、ロボティクス、AI、信号処理をまたがる領域として提案

■ 主な活動

- 人工知能学会 AI-Challenge 研究会
- 日本ロボット学会 学術講演会のオーガナイズド・セッション
- IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) のオーガナイズド・セッション
- オープンソースソフトウェア *HARK* 講習会・ハッカソン



IROSロボット聴覚セッション



日本ロボット学会オーガナイズドセッション



HARK講習会

7

ロボット聴覚デモ (同時発話認識) [Nakadai 2013]

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).

9

■ HRI-JP Audition for Robots with Kyoto University

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



hark = listen (Old English)

研究用途は無償
(商用はライセンス対応)

■ 2008年4月より、以下を目的として公開開始

- ロボット聴覚分野の活性化
- 分野間融合研究のためのツール
- ユーザからのフィードバックによる高性能化・安定性向上

■ 毎年1回ペースのリリースと無料講習会、ハッカソンの開催

Dec., 2018 : 3.0 release 予定

- ・ アーキテクチャー新
- ・ 第15回講習会: 2018/12/4 早稲田大学
- ・ 第5回ハッカソン: 2018/12/5 早稲田大学



HARK History and Tutorials

■ Apr., 2008: 初リリース (0.1.7)

- 第1回講習会: 2008/11/17 京都大学,
- 第2回講習会: 2008/12/5 韓国ソウルKIST

■ Nov., 2009: 1.0.0 プレリリース

- 第3回講習会: 2009/1/17 慶應義塾大学日吉,
- 第4回講習会: 2009/12/7 仏パリUPMC

■ Nov., 2010: 1.0.0 release : 音源分離の高性能化, ドキュメント充実

- 第5回講習会: 2010/11/25 京都大学

■ Feb., 2012: 1.1 release : 音源分離の高性能化, 64bit 対応, ROS 対応

- 第6回講習会: 2012/2/29 仏パリUPMC,
- 第7回講習会: 2012/3/9 名古屋大学

■ Mar., 2013 : 1.2 release : Kinect, PSEye 対応

- 第8回講習会: 2013/3/19 京都大学

■ Oct., 2013 : 1.9.9 release : Windows & HarkDesigner α 版

- 第9回講習会: 2013/10/2 仏ツールーズCNRS-LAAS

■ Dec., 2013 : 2.0 release : Windows & HarkDesigner 対応

- 第10回講習会: 2013/12/5 早稲田大学

■ Nov., 2014 : 2.1 release : 自己雑音抑圧対応

- 第11回講習会: 2014/11/19 早稲田大学, 11/20 第1回ハッカソン

■ Nov., 2015 : 2.2 release : 音再生モジュール追加等

- 第12回講習会: 2015/11/10 早稲田大学, 11/11 第2回ハッカソン

■ Dec., 2016 : 2.3 release : Kaldi サポート

- 第13回講習会: 2016/12/6 早稲田大学, 12/7 第3回ハッカソン

■ Dec., 2017 : 2.4 release : MVDR サポート

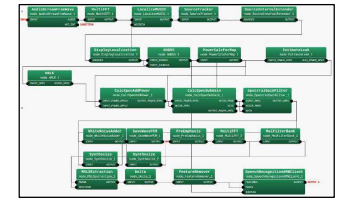
- 第14回講習会: 2017/12/5 早稲田大学, 12/6 第4回ハッカソン



HARKの特長

誰でも使える (ユーザフレンドリ)

- GUIプログラミング機能
- 容易なインストール, ドキュメントの充実
- Python サポート (HARK-Python)
- Julius, Kaldi (DNN-ASR) との連携



HARK-Designer
Chrome / Safari / Firefox
on Linux / Windows / Mac

オンライン実時間処理

- マルチチャネル A/Dを用いたオンライン処理
- ROSとのシームレスな統合
- 廉価版マイクロホンアレイのリリース

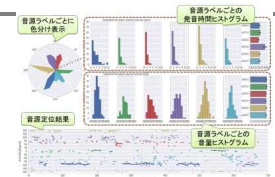


8ch アレイ Tamago (29,800円)

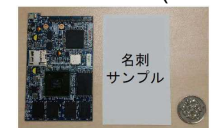
11

HARKの展開

- HARKクラウドサービス化 (HARK-SaaS)
- 組込み用ボード (RASP-MX, 販売中)
- タブレット (聴覚障がい支援, 多言語会話支援)
- 車載応用 (トークボタンレス IVI)
- 動物行動学への応用 (カエル・鳥の歌解析)
- レスキューロボットへの応用



HARK-SaaS (水本+'15)



組込みボード (中臺+'15)

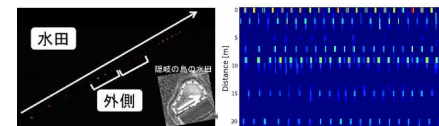


タブレット応用 (中臺+'15)

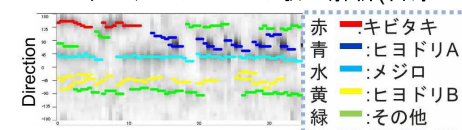


IVI 応用 (中臺+'15)

13



ニホンアマガエルの歌の解析 (合原+'14-)



鳥の歌の解析 (小島+'15, 松林+'15)

レスキューロボット：背景



Sichuan, China
'08/05



Great East Japan
Earthquake '11/03



Kumamoto, Japan
'16/04

■ 地震など世界中で多くの災害が発生

- 道路寸断 → 緊急車両通行困難

■ アクション: "The Faster, The Better"

- 72時間以内に救助が必要
 - **Golden 72 hours** (日本)
 - **The rule of threes** (欧米)
- 「アクションが一日早ければ、復旧は6か月早くなる」(Prof. R. Murphy)

■ JST ImPACT タフロボティクスチャレンジ (2014-2019)

14

ドローン聴覚の課題

1. マイクロホンアレイ収録音からの音源探索

- 音源定位・検出



2. 探索した音源のうち、音声、もしくは人に纏わる音源の判別

- 音源識別／同定



極限音響環境下で、「いつ」、「どこ」、「なに」を抽出

16

ロボット聴覚からドローン聴覚へ

■ マイクロホンアレイ搭載 UAV (Unmanned Aerial Vehicle) による災害現場での人々の搜索

- Unmanned Aerial Vehicle (UAV)

- ・ 道路が寸断されていても高速に広範囲を移動できる。

- Microphone array

- ・ カメラではとらえられない隠れた人や人がいる証拠を見つけることができる(カメラと併用)



UAV with a microphone array



関連研究 (音源探索)

■ Acoustic Vector Sensor (AVS) [Kaushik+ 05]

- 小型, 高価
- 主に軍事用途で使用(戦車、戦闘機の発見)



■ Microphone Array [奥谷+11, Basiri+ 12, Okutani+ 12, Furukawa+ 13, Ohata+ 14]

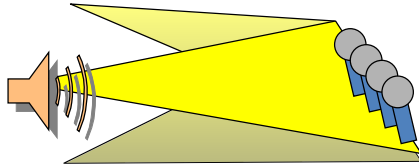
- UAV搭載マイクロホンアレイは、2011年から精力的に研究
- マイクロホンアレイ処理については、オープンソースのソフトウェアも複数あり、構築は比較的容易
 - **HARK** (Honda Research Institute Japan, Co. Ltd. Audition for Robots with Kyoto Univ.) [Nakadai 08]



Quadrotor with 16 mics

どうやって音を見つけるか？

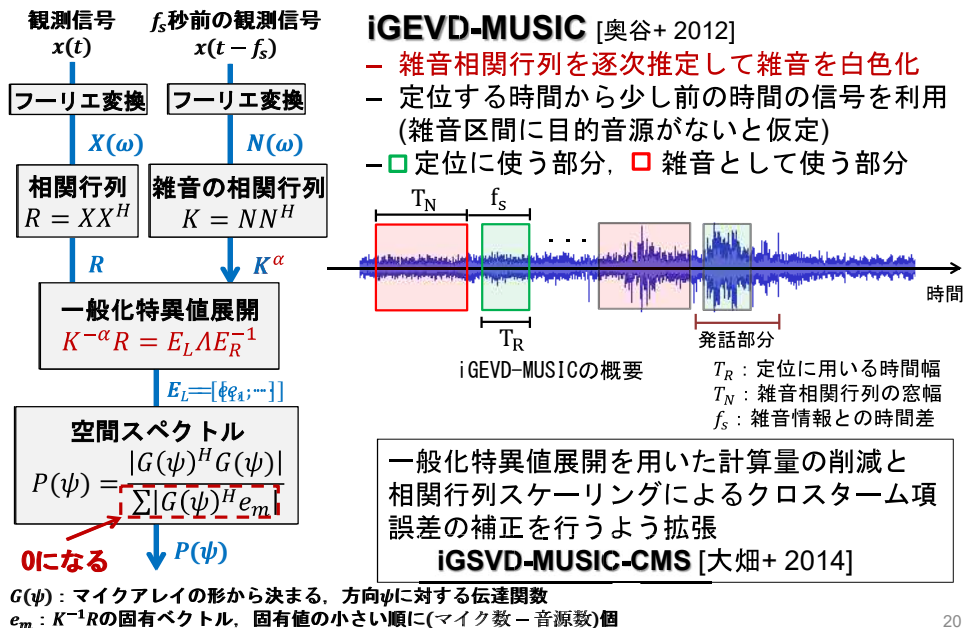
■ マイクロホンアレイ処理:



- 両耳長処理 (Jeffress モデル)
- 相互相関 (GCC-PHAT, CSP, ...)
- ビームフォーミング (WS-BF, MVDR, ...)
- MUSIC

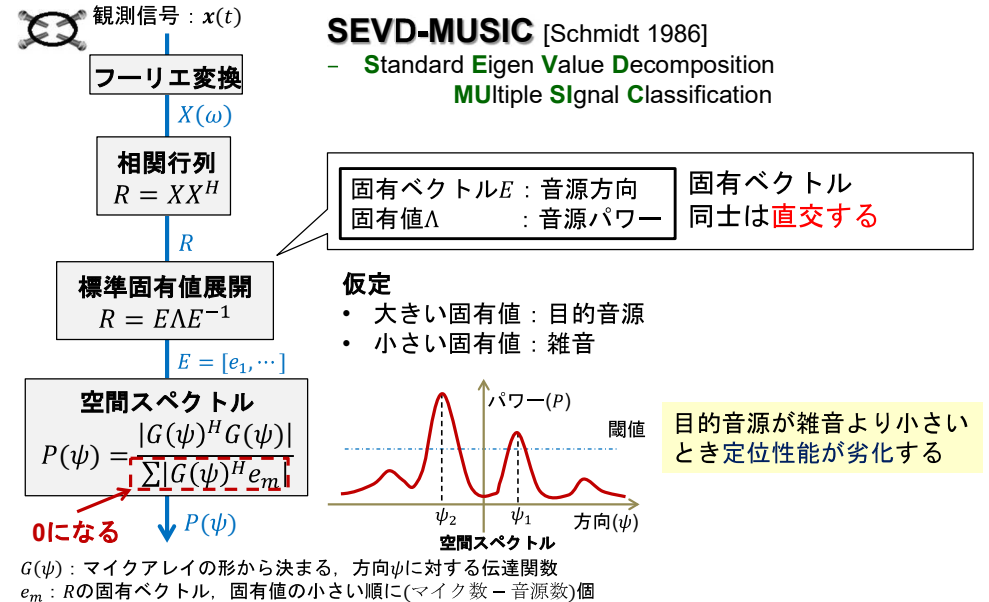
18

MUSICの拡張



20

マイクロホンアレイによる音源定位: MUSIC法



19

音源定位・検出の性能

実験環境

- ・ 場所: 屋外
- ・ 機体: 位置を固定, プロペラは回転
- ・ 音源-クアドロコプタ間距離: 4-20[m] (1[m]間隔)
- ・ 音源: 三種類音源 (音声, 車のクラクション, ホイッスル)



クアドロコプタ



音源定位手法

SEVD-MUSIC
 iGSVD-MUSIC
 iGSVD-MUSIC-CMS ($\alpha = 0.5$)

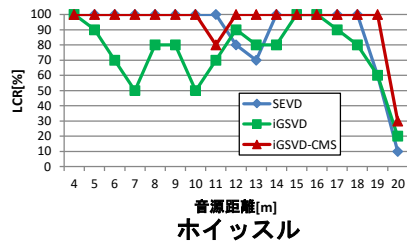
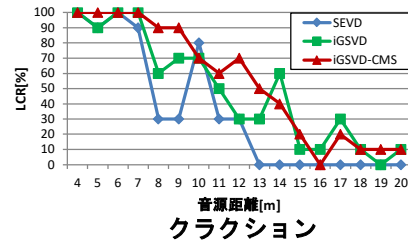
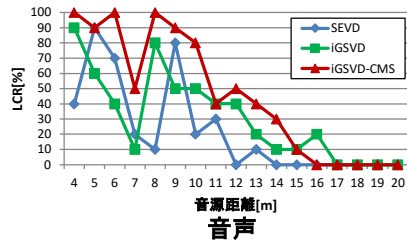
評価指標

$$LCR = \frac{\text{定位成功率}}{\text{リファレンス}}$$

目的音源をどれだけ取りもらさず定位できるか

21

音源検出距離の評価



結果

- 近距離(高SN比)の定位性能は高い
- 検出可能距離(LCR>50[%])
 - 音声12~13m
 - クラクション12~13m
 - ホイッスル18m



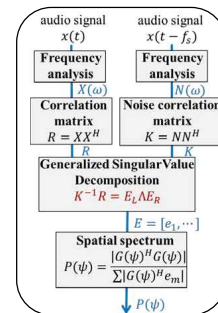
12~13m 程度の距離であれば検出可能

22

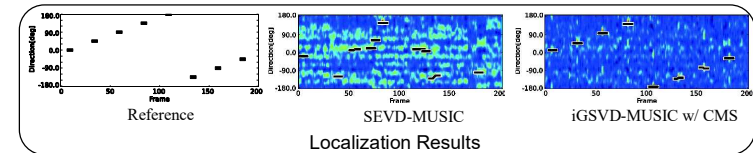
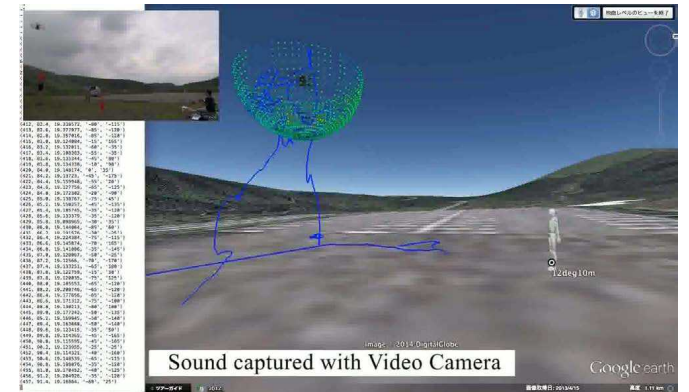
音源定位: iGEVD-MUSIC-CMS 法 [長峰+ 2014]



Quadrator with 16 mics



SSL with iGEVD-MUSIC



-15dB程度の信号でも定位可能

23

オンラインデモに向けた対策

1. 組み込み版マイクロホンアレイ処理の開発

- 球形デバイス内にARMベースの処理ボードに組み込み版HARKを搭載

→ リアルタイム処理化、
データ伝送量低減 (従来の1/100以下)

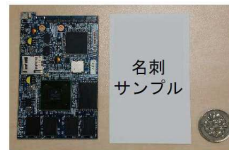
2. 三次元音源定位手法の開発

- 確率的手法を導入し、2次元の定位結果から3次元位置を推定

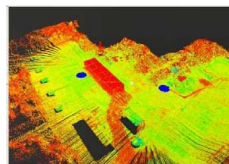
→ 地図上に音源位置を可視化

3. ケーブル一本で接続可能な、耐水仕様マイクロホンアレイデバイス開発

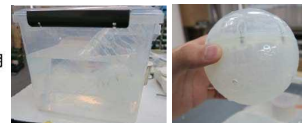
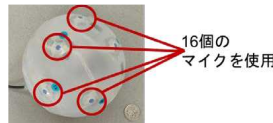
→ 雨天時にも使用可能



RASP-MX
(システムインフロンティア社)

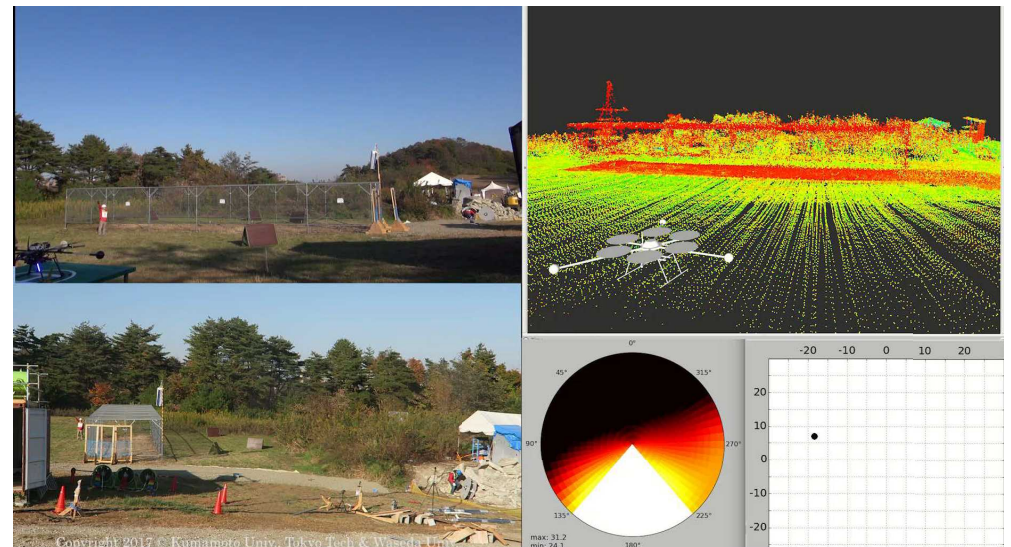


地図 (ポイントクラウド) 上に
音源位置 (青丸) 実時間表示



耐水試験
(12時間の水没試験にパス)

オンラインデモ(2017/11/11 TRC評価会)

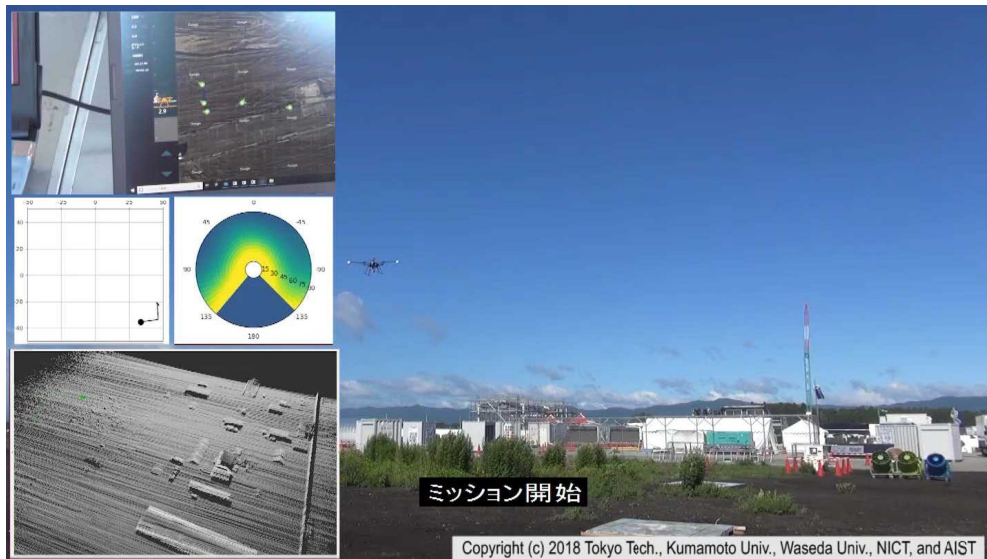


瓦礫の内外にいる人の音声を検出し、位置を地図上に表示

25

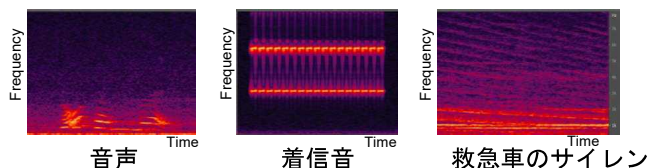
さらなる展開 (2018/6/14 TRC評価会)

- 見通し外の音源の探索
 - 極限通信チームの技術との統合
 - 中継ドローンを経由して約1km離れたところからオペレーション



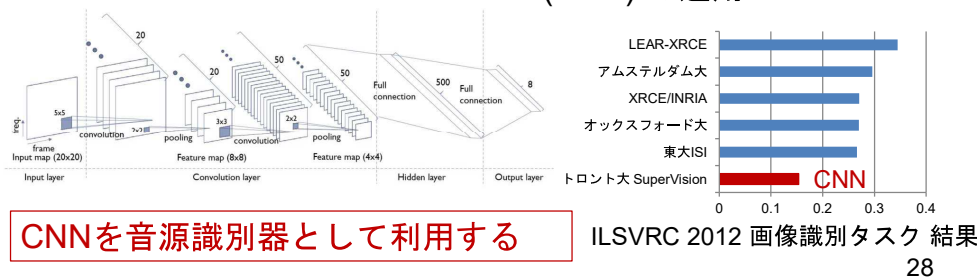
深層学習ベースの音源同定

- 音源の特徴はサウンドスペクトログラム上のパターンとして現れる。



- 音源識別は画像識別問題と捉えることができる

– Convolutional Neural Network (CNN) の適用



CNNを音源識別器として利用する

課題

1. マイクロホンアレイ収録音からの音源探索

– 音源定位・検出



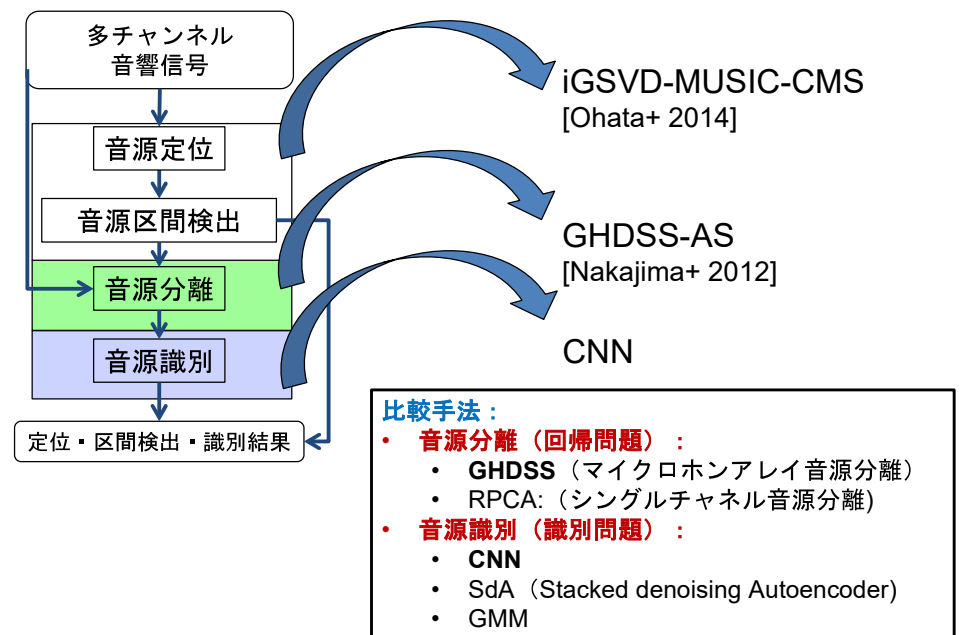
2. 探索した音源のうち、音声、もしくは人に纏わる音源の判別

– 音源識別／同定



極限音響環境下で、「いつ」、「どこ」、「なに」を抽出する。

音源検出・音源同定のアーキテクチャ[上村2015]

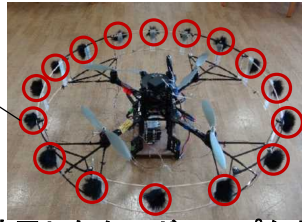


実験環境・評価指標

実験環境:

- ・ 収録に用いた機体
 - クアッドコプタ
 - ・ 16 ch マイクロホンアレイで収録
 - 機体の位置は固定、プロペラは回転
- ・ 音源位置
 - 屋外・クアッドコプタから3m離れて収録
- ・ 音源種類(10種類)
 - 音声 2 種類
 - 非音声 8 種類
- ・ 音源のSN比
 - -5dB 程度

マイクロホン



使用したクアッドコプタと
マイクロホンアレイ

学習条件

パラメータ	SdA		CNN
	pre-training	fine-tuning	
学習率	0.10	0.10	0.10
サンプル数	25,884*	25,884*	25,884*
学習回数	500	500	500
Batch size	100	100	100

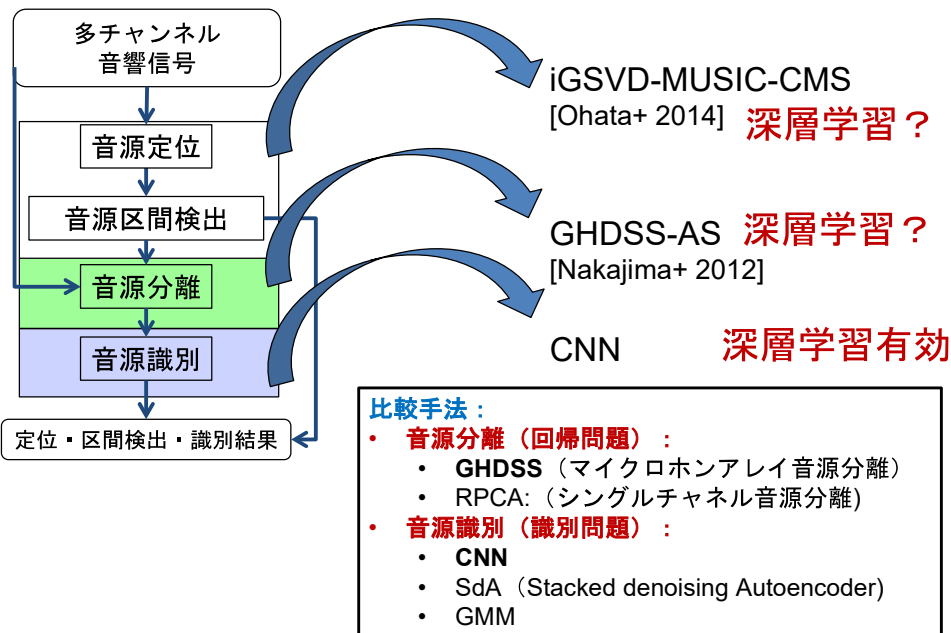
*5分割交差検証で評価

評価指標 (単体識別性能) :

(Frame Identification Correct Rate)

$$FSCR = \frac{\text{識別成功フレーム数}}{\text{定位検出フレーム数}}$$

深層学習へのシフト



実験結果

音源分離・識別手法ごとの識別率評価(%)

		音源分離手法	
		GHDSS	RPCA
音源識別手法	GMM	66.2	50.4
	SdA	58.7	40.8
	CNN	81.5	44.7

GHDSS による音源分離とCNNによる音源識別の組み合わせが**良い識別率(FSCR)**を示した

31

深層学習による音源定位 [Nelson+ 2016]

・ Deep Residual Network [He+ 2015]

- ・ ILSVRC 2015 の Winner
- ・ CNNの一種
- ・ 残差を学習することによって容易に最適化, 多層化

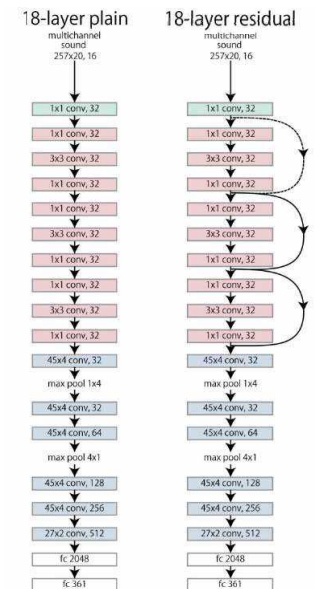
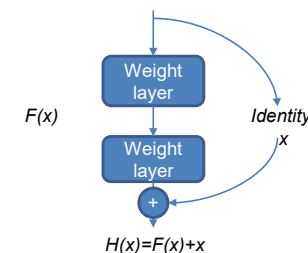
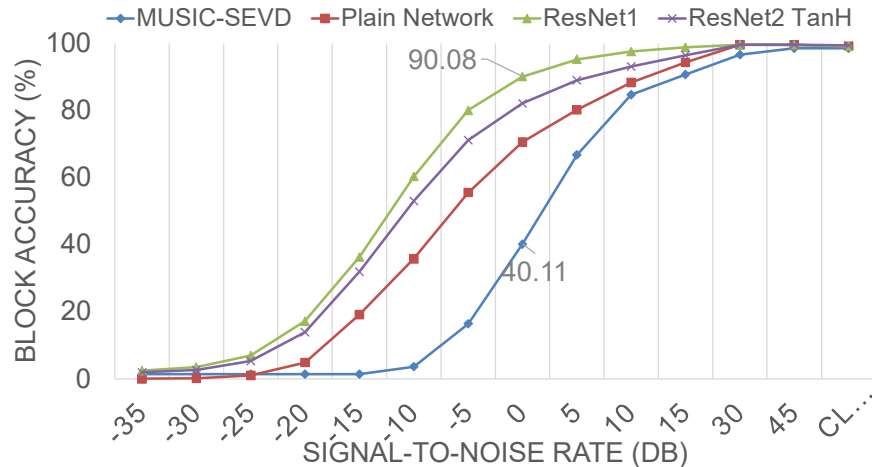


Fig. 4 Network Architecture
Left: Plain Network Right: Residual Network. The dotted line shortcut represents a 1x1 convolutional layer 33

音源定位性能



- 学習データ: JNAS (3,350,000 ファイル),
 - ロボットのモータノイズを重畳(SNR clean~-35dB)
 - 入力: STFT 257 dim
 - 出力: 36dim (定位分解能 10度)
- テストデータ: 7,200 ファイル (一方向あたり 200 ファイル)

音源定位も
深層学習？

34

音源分離性能 (音声認識)

	男性	女性	平均
処理なし	19.15	27.41	23.28
GHDSS+HRLE (HARK)	60.64	67.41	64.03
深層学習 (Denoising)	58.81	68.09	63.45

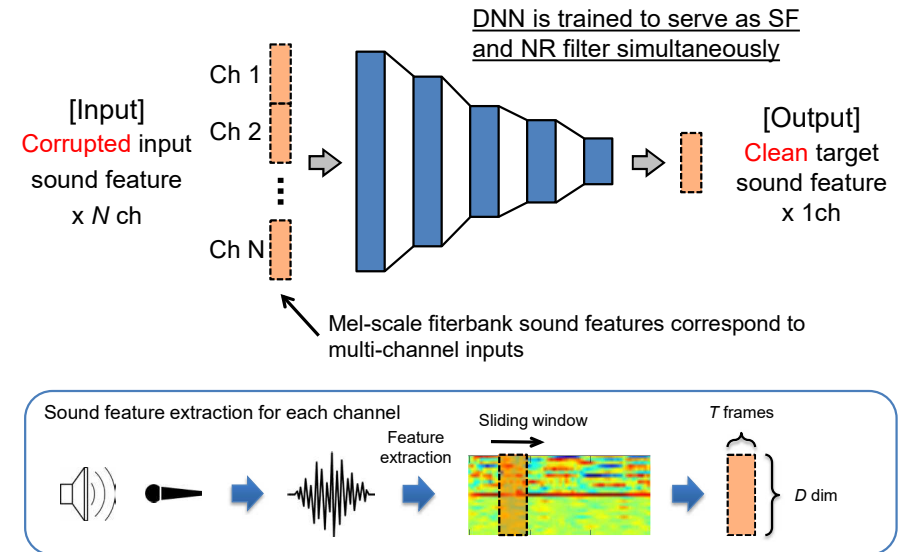


音源分離も
深層学習？

- 環境
 - 残響: 200 msec
 - 部屋のサイズ: 4 x 7 m
 - マイク数: ロボット頭部搭載 8 ch
 - 音声認識: Julius
- 評価セット
 - JNAS 200文, 二話者同時発話
 - SNR: 0 dB, 6 dB, 12 dB

36

深層学習による音源分離 [Noda+ 2015]

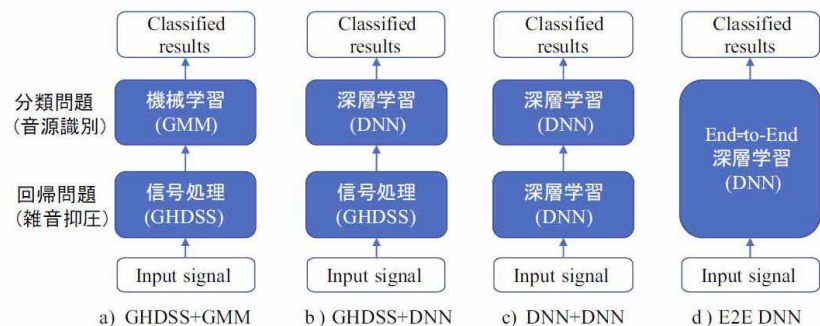


深層学習の課題

- 信号処理と深層学習の組合せ vs End-to-End 学習
 - Deep Speech, Deep Speech 2 (Hannun 2014, Amodei 2015)
 - WaveNet (Google DeepMind 2016)
- 大量の教師データが必要
 - アノテーションは人手
- 処理時間がかかる
 - 学習はオフラインでもいいが、識別／回帰はオンラインがいい

37

End-to-End がよいのか？ [中臺18 ASJ]



手法	識別性能
a) GHDSS+GMM	95.90%
b) GHDSS+DNN	99.34%
c) DNN+DNN	97.49%
d) End-to-End DNN	96.83%

Parrot Bebop Drone w/ 8ch mic-array
#samples:60k, 5-fold CV, input:1,600-dim
SNR: 0dB, batch:100, epoch:300, 6 layered DNN



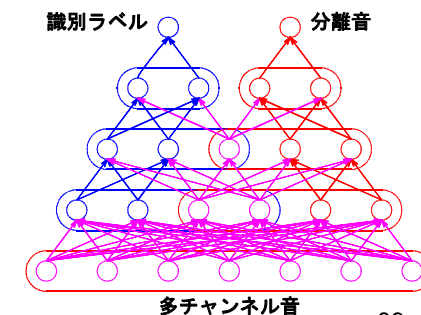
深層学習：
パワフル
実応用考慮：
安直な E2E より
音響信号処理と
深層学習の
いい所どりを

38

Partially Shared Deep Neural Network (PS-DNN) [森戸+ 2016]

- Multi-purpose DNN の一種
 - 二つのDNNの中間層を一部共有するように結合
- 音源識別 (classification)
 - 多チャンネル信号を入力として、音源名の識別学習
→ End-to-End モデル
- 音源分離 (regression)
 - 多チャンネル信号を入力として、デノイズング
 - 正解のクリーンデータがなくても
信号処理の音源分離結果を代用可
(手動アノテーション不要)

- 共有部分は音源名のアノテーション (教師データ) が無くても学習可能
- 音源識別器の学習が少量のアノテーションで学習できる

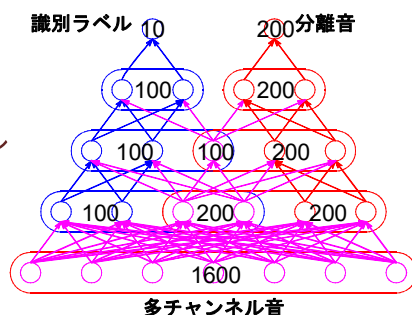


39

PSA の効果

■ PS-DNNの構成

- 入力: 音響特徴量20次元 x 10 フレーム
x 8チャンネル = 1600次元
- 出力: 音響特徴量200次元 x 1チャンネル
識別ラベル10次元



- DNN≡PS-DNN-
音源分離学習は音源識別学習を阻害しない
- PS-DNN- < PS-DNN
未アノテーションデータを効果的に利用

学習データ量/ アノテーション 有データ	DNN	PS-DNN-	PS-DNN
100%	98.87%	98.93%	-
80%	97.36%	97.11%	97.61%
60%	94.81%	94.95%	96.01%
40%	90.76%	91.34%	93.00%
20%	86.60%	86.96%	87.83%

DNN: 共有なし, end-to-end 識別学習
PS-DNN-: アノテーション有データのみ使用
PS-DNN: 全データ使用

40

デコードの高速化 [Takeda+17]

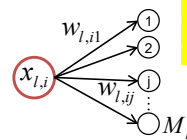
- DLベースの手法はデコードにも計算時間がかかる(組み込み用のプロセッサでも動作させたい)
- アプローチ: エントロピーベースの重みの量子化とノードの枝刈り

重みエントロピー

$$q(x_{l,i}) = -M_{l,i} \sum_w p_i(w) \log p_i(w)$$

$$p_i(w) : \text{正規化ヒストグラム}$$

$$w_{l,i,j} \quad (j = 1, \dots, N_{l,i})$$



エントロピーが大きい
- コネクション数が多い
- 量子化ビット数が多い

$w_{l,i,j}$: n-bit integer
e.g. 1-bit [0, 1]
2-bit [00, 01, 10, 11]

ノードエントロピー

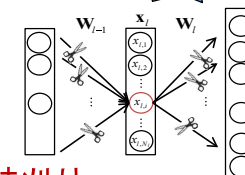
$$q(l, i | D) = -p_{0,l,i} \log p_{0,l,i} - p_{1,l,i} \log p_{1,l,i}$$

$$p_{0,l,i} = N_{0,l,i} / (N_{0,l,i} + N_{1,l,i})$$

$$p_{1,l,i} = N_{1,l,i} / (N_{0,l,i} + N_{1,l,i})$$

- $N_{0,l,i}$: シグモイド出力が 0.5 以下のサンプル数
- $N_{1,l,i}$: シグモイド出力が 0.5 以上のサンプル数

エントロピーが大きい
ノードは貢献が大きい



エントロピーに基づいたノードの枝刈り

枝刈りの効果

Pruned nodes (%)	# of param. (M)	10ビット量子化固定				2ビット, 8ビット混合
		ノード エントロピーのみ		ノード&重み エントロピーのみ (提案法)		RTF
		RTF	WA (%)	RTF	WA (%)	
0	9.71	0.941	81.86	0.941	81.86	0.941
30	4.98	0.428	82.02	0.390	81.89	0.204
40	4.00	0.344	81.72	0.326	81.53	0.188
50	3.17	0.273	80.74	0.254	80.68	0.171
60	2.48	0.213	79.34	0.197	79.29	0.154

- 音声認識タスク (Juliusベース単語認識)
- 30%のノード枝刈りを行った際に、認識性能を維持したまま提案法では実時間性を0.428→0.390まで向上
- 量子化ビット数を工夫すれば、さらに0.204まで実時間性が向上

42

お知らせ

- Oct. 5 スペイン, マドリード
 - HARK Tutorial @ IEEE/RSJ IROS 2018
- Dec. 3 早稲田大学 (西早稲田キャンパス)
 - AIチャレンジ研究会 (ロボット聴覚特集号)
- Dec. 4 早稲田大学 (西早稲田キャンパス)
 - HARK 講習会
- Dec. 5 早稲田大学 (西早稲田キャンパス)
 - HARK ハッカソン

44

まとめ

- ロボット聴覚の紹介
- ロボット聴覚からドローン聴覚への展開
 - 信号処理ベースの音源探索
 - 深層学習ベースの音源同定
 - 深層学習の課題と対策

43