

# 複数ロボットによる音源定位結果を統合し発話者を特定するシステム

## Speaker Identification System

### Based on Sound Source Localization Results from Multiple Robots

中島大一, 駒谷和範, 佐藤理史

Taichi Nakashima, Kazunori Komatani, Satoshi Sato

名古屋大学大学院 工学研究科

Graduate School of Engineering, Nagoya University

{taichi\_n, komatani, ssato}@nuee.nagoya-u.ac.jp

## Abstract

Humanoid robots need to head toward human participants when answering to their questions in multi-party dialogues. Some positions of participants are difficult to localize from robots in multi-party situations, when the robots can only use their own sensors. We present a method of identifying who is speaking more accurately by integrating the multiple sound source localization results obtained from two robots. This method employs two robots and places them so as to compensate for each other's localization capabilities and then integrate their two results. Our experimental evaluation revealed that using two robots improved speaker identification compared with when only one robot was used. We furthermore implemented our method using humanoid robots and constructed a demo system.

## 1 はじめに

ヒューマノイドロボットを用いた複数人会話システムの開発を行っている。複数人会話システムとは、2人以上のユーザと会話するシステムである[12]。

今日までに開発されてきた複数人会話システムには2つの問題がある。1つ目の問題は、多くのシステムが特殊なデバイスを利用しており、特定の環境でしか動作しないことである。複数人会話システムは、ユーザが何を発話したのかだけでなく、そのユーザがどこにいるのか、そして誰に対して発話したのかを特定する必要がある。このようなユーザの行動を検出、追跡するために、超音波センサ[7]や、高解像度 (HD) カメラ[4]、広角カメラ[2]といったデバイスが利用されている。もしくは、カメラやマイクロ

フォンが多く設置されたスマートルーム[5]が前提とされてきた。2つ目の問題は、ユーザがシステムに対して何を発話すればよいのかが分からず、会話が中断してしまうことである。この問題は、単一のユーザを相手にする対話システムでも発生する[3]ため、複数人会話システムでも起こりうる問題である。

本研究では以下のアプローチでこれらの問題の解決を図る。

1. ロボットに搭載されたマイクとカメラのみを利用する。つまり、特殊なデバイスが準備された特別な環境を仮定しない。
2. 2体のロボットを利用する。ユーザが発話を止めたときに、ロボット同士で会話を行うことで会話を持続させる。さらに、ロボットに搭載されたカメラやマイクロフォンの能力は十分ではないため、2体を用いることでそれを補う。

複数人会話を行う状況として、図1に示すような複数ユーザが机を囲み、その机の上に配置されたロボットと行う会話を想定する。この状況設定は、複数人会話システムにおけるユーザの位置の決定を単純化するものである。つまり、複数のユーザが机を囲んだ状況では、ユーザの位置を机の周辺に限定できる。

本論文では、複数人会話システムの実現の第一歩として、発話したユーザがどこにいるかを特定する「発話者の特定」を扱う。発話者を特定できれば、ロボットに発話したユーザに対して顔を向け、応答を行わせることができる。この挙動は、ユーザに自分が聞き手であることを自覚させ[8]、かつ会話に参加していると感じさせる[1]ことが期待できるため、複数人会話システムにとって重要な要素である。このようなロボットの挙動に加え、複数人会話ではシステムが個々のユーザに質問したり、話題を提供したりするのも重要である。例えば、まだ発話をしていないユーザに対して発話を促せるのが望ましい。このために

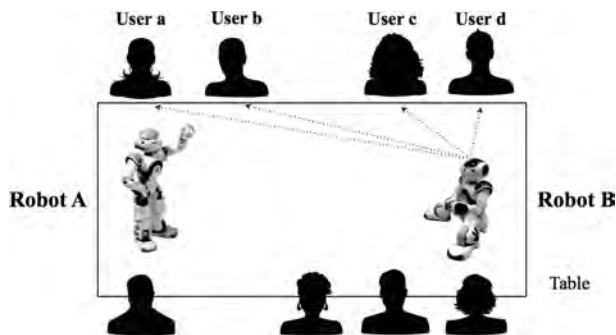


Figure 1: 想定する会話状況。複数ユーザが机を囲み、机の上に配置したロボットと会話を行う。

は、システムは個々のユーザの位置を個別に特定する必要がある。

本研究では音源定位を用いて発話者の特定を行う。音源定位とは、音源がどの方向から到来したかを推定する手法である。ここで、図1に示す本稿で想定する会話状況では、以下の2つの問題がある。

1. ロボットに搭載されたマイクのみを用いる場合、定位が困難な位置が存在する。
2. 雑音の誤検出が避けられない。つまり、音源定位結果が常に発話者の位置を示しているとは限らない。

本稿で想定する会話状況において、ロボットから遠く離れるほど、2人のユーザの位置を示す角度差は小さくなる。例えば、図1において、Robot Bから遠い位置に存在するUser aとUser bの間の角度差は小さい。このため、2人のユーザを個別に定位するのは難しい。

本研究では、2体のロボットを図1のように対面的に机に配置し、それらから得られる音源定位結果を統合することで、これら2つの問題の解決を図る。この配置により、2体のロボットのマイクロフォンによる音源定位能力を互いに補うことができる。例えば、Robot Bの位置からは、User aとUser bの間に十分な角度差はないが、Robot Aの位置からは十分な角度差があり、それらを個別に定位することは容易である。この2体のロボットから得られる音源定位結果に対して、パワー（音圧）で重みを付け、統合を行う。この音源定位結果の正しさを示す指標としてパワーを用いる。

さらに、本統合手法をヒューマノイドロボットNAO<sup>1</sup>を用いて実装し、デモシステムを構築した。本システムは発話したユーザを特定し、音声認識に基づきそのユーザに対して顔を向けて応答を返す。統合された定位結果のパワーが低い場合は、ロボットに搭載されたカメラを用いて発話者の存在を確認する。

<sup>1</sup><http://www.aldebaran-robotics.com/en/>

## 2 関連研究

複数人会話において発話者の特定を行う最も単純な手法は、ユーザそれぞれにマイクロフォンを持たせることである[6]。マイクロフォンとユーザの位置からシステムは発話者がどこにいて、いつ発話したのかを特定できる。この手法は会話を開始する前に、各ユーザにマイクを準備する必要がある。他にも、多くのマイクロフォンやカメラを設置したスマートルームで会話を行うことで、発話者を特定する手法も考えられる[11]。この手法では、室内に設置された複数のセンサを利用して、ユーザの行動を追跡し、発話者を特定する。この手法は、このような特定の環境でのみ有効である。

これに対して、本研究ではロボットに搭載されたマイクロフォンやカメラのみを用いて発話者の特定を行う。この場合、本稿の想定する会話状況では以下の2つが前提となる。

1. ユーザがカメラの視野角内に常に存在するとは限らない。これはロボットに搭載されたカメラの視野角は狭く、解像度も低いためである。
2. 狭い視野角を補うために、常に周りを見回し続けるのは不適當である。これはロボットが発話の当事者であり、このような挙動は会話として不自然になるためである。

Faihらは、唇領域を検出することで、発話者を特定する手法を提案した[4]。Bohusらのシステムは、画像情報に基づいてユーザを追跡し、ジェスチャなどのユーザ行動を認識することで発話者の特定を行う[2]。これらの手法はユーザが常にカメラの視野角内に存在する場合には有効である。Bennewitzらは、ユーザがロボットのカメラの視野角外に存在する場合であっても、顔検出に基づき参加者の存在について確率的な信頼度を維持する手法を提案した[1]。この手法では、常にユーザをカメラの視野角内に捉えておく必要はないが、一度もカメラの視野角に入らない発話者の位置を特定するのは困難である。

本研究では、主に音源定位結果に基づき発話者の特定を行う。音源定位結果を用いれば、ロボットのカメラの狭い視野角内にユーザが存在することが仮定できなくても、発話者の特定ができる。さらに、一度も視野角内に入らない発話者の特定も期待できる。

## 3 音源定位結果の統合

本研究では2体のロボットから得られる音源定位結果を統合する。本章ではその統合方法について述べる。

2体のロボットのマイクロフォンを通して、音源の到来方向を角度で示す音源定位結果とそのパワー（音圧）を得る。音源定位結果は、ロボットの正面方向を0度とし、反時計回りを正方向として得られる（例えば、ロボットの左

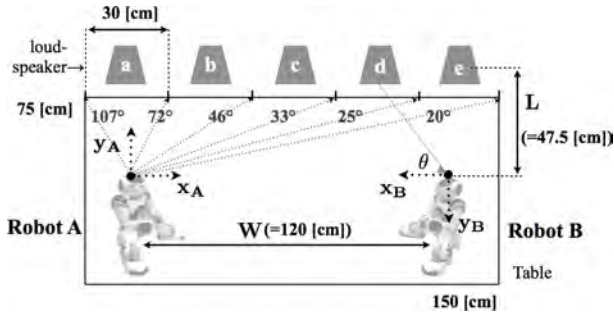


Figure 2: ロボットとスピーカの配置. スピーカをユーザが存在すると考えられる位置に配置する.

手方向は90度である). この設定で, 以下の3つのステップにより2体のロボットから得られる音源定位結果を統合する.

1. 2体のロボットで異なる座標系を一致させる.
2. 音源定位結果に対して, パワーで重みを付ける.
3. 2体のロボットから得られた重み付きの音源定位結果を足し合わせる.

まず, 図2に示す2体のロボットの座標系を一致させる. 図2における各変数を以下のように定義する.

- $(x_R, y_R)$ : ロボットごとの座標. ここで  $R \in \{A, B\}$  とし,  $A$  と  $B$  は図2の Robot A, Robot B と対応する. 座標の原点はロボットの頭部である.  $x_R$  軸の正方向はロボットの正面方向とし,  $y_R$  軸の正方向は正面から見て左方向である. 本研究ではロボットは同じの高さの平面上に存在すると仮定し, 垂直方向については考慮しない.
- $W$ : 2体のロボット間の水平距離.
- $L$ : ロボットとユーザが存在すると考えられる位置の最短の水平距離.

ここでは, Robot B の座標系を Robot A の座標系に一致させる. Robot B が音源定位結果  $\theta$  を得たとき, Robot B の座標系における音源の座標  $(x_B, y_B)$  は式1で表せる.

$$(x_B, y_B) = \left( \frac{l}{\tan \theta}, l \right) \quad (1)$$

$$l = \begin{cases} L & , \text{ if } \theta > 0 \\ -L & , \text{ if } \theta < 0 \end{cases}$$

Robot A の座標系における  $(x_B, y_B)$  は式2で表せる. ここで,  $\alpha$  は2体のロボットの座標系間の回転角度差であり,  $(u, v)$  は, Robot A の座標系における Robot B の原点座標を示す.

$$\begin{pmatrix} x_A \\ y_A \end{pmatrix} = \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix} \begin{pmatrix} x_B \\ y_B \end{pmatrix} + \begin{pmatrix} u \\ v \end{pmatrix} \quad (2)$$

これにより, ロボット間で異なる座標系の一致をとる. 実際に, ロボットを図2のように配置したとき, 各変数の値は

それぞれ,  $L = 48.5$  [cm],  $\alpha = 180$  [度],  $(u, v) = (W, 0)$ , そして  $W = 120$  [cm] であった. 現在,  $L$  や  $\alpha$ ,  $W$  は人手で計測している.

次に, 音源定位結果に対して, パワーで重みを付ける. 雑音による音源定位結果のパワーは小さいと仮定し, パワーを音源定位結果の正しさを示す指標として利用する. 音源定位結果  $\theta_r$ , パワー  $p_r$  が得られた場合, 音源定位結果の曖昧さは正規分布に従うと仮定し, 確率密度関数  $f_r(\theta)$  を定義する (式3). ここで,  $r$  は ID を示す, 例えば, Robot A の音源定位結果は  $\theta_A$  と表現する. 式3において,  $\sigma_r^2$  は分散であり, 音源定位結果がどれだけ不確かであるかを示す.

$$f_r(\theta) = \frac{1}{\sqrt{2\pi\sigma_r^2}} \exp\left(-\frac{(\theta - \theta_r)^2}{2\sigma_r^2}\right) \quad (3)$$

確率密度関数の最大値  $f_r(\theta_r)$  は音源定位結果  $\theta_r$  のパワー  $p_r$  に比例すると仮定する (式4). この仮定は, パワー  $p_r$  が大きいほど, 音源定位結果が  $\theta_r$  である確率が大きくなることを示す. ここで, 式4の  $C$  は定数であり実験的に決定する.

$$f_r(\theta_r) = \frac{1}{\sqrt{2\pi\sigma_r^2}} = \frac{1}{C} p_r \quad (4)$$

式4より  $\sigma_r$  を定める (式5). 式5より,  $\sigma_r$  はパワー  $p_r$  に反比例する. つまり, パワーが大きいほど音源定位結果は散らばりが小さいとしている.

$$\sigma_r = \frac{C}{\sqrt{2\pi}} \frac{1}{p_r} \quad (5)$$

確率密度関数  $f_r(\theta)$  の例を図3に示す. グラフの横軸は, 音源定位結果を示し, 縦軸はその確率を示す.

最後に, 上記のステップの後に得られた音源定位結果を足し合わせる. Robot A, Robot B から音源定位結果  $\theta_A, \theta_B$  とそのパワー  $p_A, p_B$  が得られたとき, それぞれ確率密度関数  $f_A(\theta), f_B(\theta)$  を定義する.  $f_A(\theta)$  と  $f_B(\theta)$  を式6, 7に適用し, 統合による音源定位結果  $\theta_{mix}$  とそのパワー  $p_{mix}$  を得る.

$$\theta_{mix} = \arg \max_{\theta} (f_A(\theta) + f_B(\theta)) \quad (6)$$

$$p_{mix} = C(f_A(\theta_{mix}) + f_B(\theta_{mix})) \quad (7)$$

さらに, 統合によって得られるパワーに閾値を設定する. そして, 閾値より小さいパワーを持つ音源定位結果を削除する. これにより, ノイズ等に基づく誤った音源定位結果の悪影響を減らすことを期待できる. この閾値は実験により最も精度の高い値を設定する.

## 4 評価実験

2体のロボットによる音源定位結果の統合により, 1体のロボットのみを用いる場合と比較して, 発話者の特定性能が向上することを確認する. 本実験では, スピーカから音声を再生し, そのスピーカの位置を特定した.

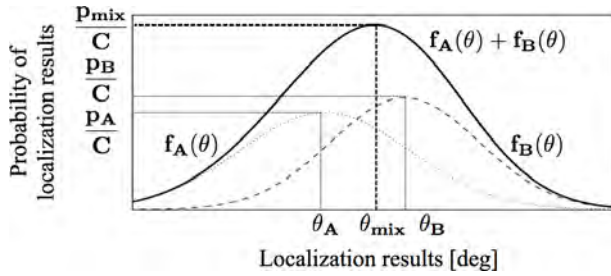


Figure 3: 確率密度関数の足し合わせの例

#### 4.1 実験設定

図2に示すように、机 (150cm × 75cm) を準備し、ユーザが存在すると考えられる位置にスピーカを配置した。スピーカはそれぞれ 30cm 間隔で配置し、その中心から ±15cm をそのスピーカの領域とした。図2には Robot A から見たスピーカの領域を示す。

本研究では、スピーカから音声を再生しているときに、音源定位結果がそのスピーカの領域内であれば、その定位結果を正解とみなす。例えば、図2において、bのスピーカから音声を再生し、Robot A で音源定位を行う場合を考える。このとき、Robot A からみてbのスピーカは 46度から 72度の間に存在するため、定位結果がその間であれば、正解とする。

音源定位には、ロボット聴覚システム HARK [9]を用いた。HARK は MUltiple SIgnal Classification (MUSIC) 法[10]に基づき、1 フレーム (0.01 秒) ごとに音源定位結果とそのパワーを出力する。MUSIC 法は、音源と入力に用いるマイクロフォン間のインパルス応答 (伝達関数) に基づき、音源を定位する。マイクロフォンには、ヒューマノイドロボット NAO の頭部の前後左右に搭載された4つのマイクロフォンを用いた。伝達関数を計算するためのインパルス応答は、マイクロフォンの中心より 1m から、10度間隔で 36 点計測した。したがって、音源定位結果の角度分解能は 10 度である。

#### 4.2 データと評価指標

スピーカから音声ファイルを再生し、それをロボットのマイクロフォンで録音したデータに対して音源定位を行った。各音声ファイルには、1名のユーザによる発話が録音されていた。発話は全部で5種類で、平均の発話長は 1.0 秒であった。それらを男女4名により録音し、図2に示す a から e の全5カ所から再生した。評価は発話ごと、つまり、全 100 発話について行った。

評価指標には Precision ( $P$ ), Recall ( $R$ ), そして  $F$  値 ( $F$ ) を用いた。本研究では、それらを以下のように定義する。

$$P = \frac{\text{発話中にスピーカの領域内を定位したフレーム数}}{\text{全検出フレーム数}}$$

$$R = \frac{\text{発話中にスピーカの領域内を定位したフレーム数}}{\text{全発話フレーム数}}$$

$$F = \left( \frac{1}{P} + \frac{1}{R} \right)^{-1}$$

#### 4.3 実験結果

5カ所 (a から e) に配置したスピーカの特定制結果を表1に示す。表は左から Robot A のみ, Robot B のみ, そして統合による特定制結果である。パワーの閾値は、Robot A のみ, Robot B のみ, 統合 ( $C = 800$ ) の各条件でそれぞれ、24, 25.5, 25 のとき、ALL の  $F$  値が最大となった。

1体のロボットのみを用いた場合には、ロボットから遠い位置のスピーカの特定制が難しい。例えば、Robot A から遠い位置にある、スピーカ c, d, e の特定制性能は低い。これは、正解とするスピーカの領域が、ロボットとスピーカが離れるほど狭くなり、その領域の定位が困難であるためである。さらに、ロボット間で特定制性能に差がある。これは、ロボットのマイクロフォンの性能が異なるためである。

統合により、1体のロボットのみでは困難であった位置の特定制が可能になっている。特に、cのスピーカはどちらのロボットも1体のみではほとんど特定制ができていないが、統合により他の位置のスピーカと同等の  $F$  値が得られている。c以外のスピーカの  $F$  値は若干の精度の低下がみられる。これは、どちらのロボットも常に正しい音源定位結果を出力するわけではなく、誤った音源定位結果が統合に悪影響を与えることがあるためである。

#### 4.4 パワーごとの音源定位結果

本研究では、統合によるパワー ( $p_{mix}$ ) を音源定位結果の正しさを示す指標として用いる。パワーごとのスピーカの特定制結果を調べることで、パワーが音源定位結果の正しさを示す指標として利用できることを確認する。図4に、パワーごとの平均誤り率とその平均角度誤差を示す。グラフの横軸は、パワーで 8dB ごとに集計している。左の縦軸は誤り率であり、右の縦軸は平均角度誤差を示す。誤り率は、音源定位結果の出力数のうち、誤りであったものの割合である。ここで、正解であるスピーカの領域以外で検出された音源定位結果を、誤りとした。平均角度誤差は、誤りであった音源定位結果と、スピーカの領域の中心との誤差の平均として計算した。

図4の平均角度誤差をみると、パワーが大きいほど平均角度誤差は小さい。つまり、パワーが大きいほど誤りであってもスピーカの領域に近い位置を定位できている。パワーが小さいときは、誤り率は高く、平均角度誤差も大きい。これより、統合後のパワーを用いることで、音源定位結果の正誤の区別が期待できる。例えば、パワーが小さい音源定位結果が得られたときは、音源定位結果は誤り

Table 1: スピーカの特特定結果

スピーカ	Robot A			Robot B			統合		
	<i>Precision</i>	<i>Recall</i>	<i>F</i> 値	<i>Precision</i>	<i>Recall</i>	<i>F</i> 値	<i>Precision</i>	<i>Recall</i>	<i>F</i> 値
a	0.56	0.89	0.69	0.00	0.00	-	0.57	0.85	0.68
b	0.49	0.65	0.56	0.00	0.00	-	0.40	0.50	0.45
c	0.00	0.00	-	0.13	0.13	<b>0.13</b>	0.38	0.49	<b>0.43</b>
d	0.06	0.03	<b>0.04</b>	0.63	0.83	0.72	0.48	0.67	0.56
e	0.09	0.03	<b>0.05</b>	0.50	0.69	0.58	0.39	0.61	0.48
ALL	0.33	0.32	0.33	0.39	0.33	0.36	0.45	0.62	0.52

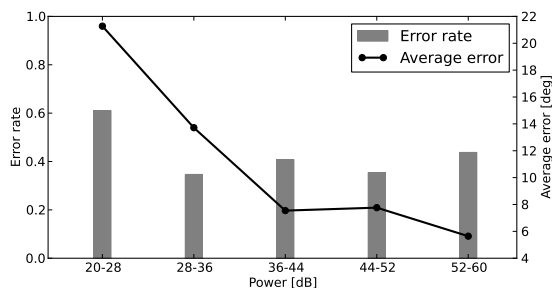


Figure 4: パワーごとの誤り率と平均角度誤差

である可能性が高い。そのため、システムは顔検出などを行うことで、発話者の存在を確認するのが望ましい。

## 5 デモシステム

本統合手法をヒューマノイドロボットを用いて実装し、デモシステムを構築した。図5に、構築したシステムと複数ユーザによるインタラクションの様子を示す。システムのタスクは我々の研究室の紹介で、ユーザはロボットに研究室に関する質問ができる(例えば、「研究室の生活について教えて」)。2体のロボットにはそれぞれ役割を設定した。役割は主にユーザの質問に答える説明役と、ユーザと共に質問を行う質問役である。構築したシステムには以下の4つの特徴がある。

- 複数ユーザの中から発話者を特定し、そのユーザに顔を向けて応答を行う(図5の上の写真)。
- 統合によるパワー ( $p_{mix}$ ) を、音源定位結果の正しさを示す指標として用いる。パワーの小さい音源定位結果が得られたとき、ロボットはその方向に向けて自身に搭載されたカメラを用いて顔検出を行い、発話者の存在を確認する。パワーが大きいときは、確認は行わず、そのまま応答を返す。
- 2体のロボットで同時に音声認識を行い、音源定位結果を用いて発話者に近いロボットが得た音声認識結果を採用する。
- 一定時間ユーザの沈黙を検出したとき、質問役のロボットが説明役のロボットに質問を行い、会話の間を繋ぐ。

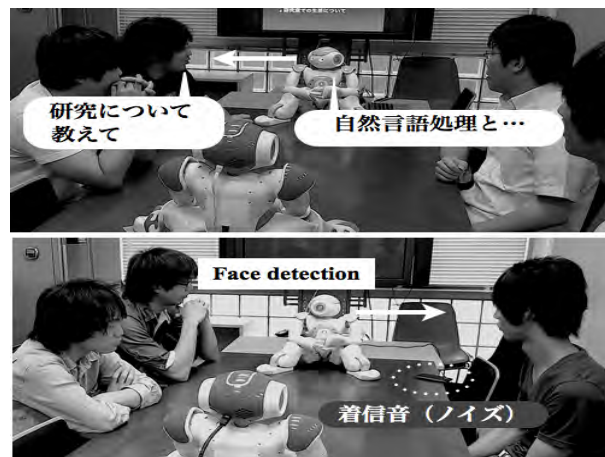


Figure 5: 複数ユーザとシステムのインタラクションの様子。上の写真では、システムが発話者を特定し、そのユーザに顔を向け応答を返す様子を示す。下の写真は、パワーの低い音源定位結果が得られた際に顔検出を行い、発話者の存在を確認の様子を示す。

構築したシステムと複数ユーザのインタラクションのデモ動画はオンラインで視聴できる<sup>2</sup>。

## 6 まとめと今後の課題

本論文では、2体のロボットから得られる音源定位結果を統合し、発話者の特定を行う手法について述べた。評価実験により、1体のロボットのみを利用する場合と比較して、統合により発話者の特定性能が相対的に向上することを示した。しかし、絶対的な発話者の特定精度は  $F$  値にして0.52程度(表1)であり、複数人会話を行うには性能向上が必要である。本研究の問題設定において、発話者の特定精度が低い理由は以下の2点である。まず第一の理由は、正解条件、つまりスピーカの正解領域を厳しく設定しているためである。例えば、図2において、Robot Aの位置からみたd, eのスピーカの間角度差は、他の位置のスピーカと比べて狭い。このような正解条件を用いるのは、複数人会話においてシステムが個々のユーザを個別に定位できることが重要であるためである。第二の理由は、マイクロフォンとスピーカが離れているためである。

<sup>2</sup>[http://sslab.nuee.nagoya-u.ac.jp/en/?page\\_id=112](http://sslab.nuee.nagoya-u.ac.jp/en/?page_id=112)

このような設定では、接話型のマイクロフォンを利用する場合と比較して、部屋の残響や環境雑音の影響が避けられない。

発話者の特定精度を向上させるための今後の課題は、発話者の存在を示す別の情報源も同時に用いることである。別の情報源には、例えば、ロボットのカメラから得られる画像情報や、会話開始からある時点までに得られた音源定位結果が考えられる。これらの情報を現在の手法と同時に用いることで発話者の特定精度のさらなる向上が期待できる。

個別のユーザの定位に基づいたインタラクションの実現も今後の課題である。例えば、まだ発話していないユーザに対して、発話を促すといった挙動を生成する。

## 謝辞

Nao と HARK を接続するプログラムは、京都大学の水本武志氏と協力して作成した。本研究の一部は、JST 戦略的創造研究推進事業さきがけの支援を受けた。

## 参考文献

- [1] Maren Bennewitz, Felix Faber, Dominik Joho, Michael Schreiber, and Sven Behnke. Integrating vision and speech for conversations with multiple persons. In *Proceedings of IEEE/RSJ the International Conference on Intelligent Robots and Systems (IROS)*, pages 2523–2528, 2005.
- [2] Dan Bohus and Eric Horvitz. Models for multiparty engagement in open-world dialog. In *Proceedings of the SIGDIAL 2009 Conference*, pages 225–234, 2009.
- [3] Alexander Gruenstein and Stephanie Seneff. Releasing a multimodal dialogue system into the wild: User support mechanisms. In *Proceedings of the 8th SIGdial Workshop on Discourse and Dialogue*, pages 111–119, 2007.
- [4] Fasih Haider and Samer Al Moubayed. Towards speaker detection using lips movements for human-machine multiparty dialogue. In *FONETIK 2012*, 2012.
- [5] Natasa Jovanovic, Rieks op den Akker, and Anton Nijholt. Addressee identification in face-to-face meetings. In *Proceedings of the 11th Conference of the EACL*, 2006.
- [6] Yoichi Matsuyama, Hikaru Taniyama, Shinya Fujie, and Tetsunori Kobayashi. Framework of communication activation robot participating in multiparty conversation. In *Proceedings of AAAI Fall Symposium, Dialog with Robots*, pages 68–73, 2010.
- [7] Samer Al Moubayed, Jonas Beskow, Mats Blomberg, Björn Granström, Joakim Gustafson, Nicole Mirnig, and Gabriel Skantze. Talking with furhat - multi-party interaction with a back-projected robot head. In *FONETIK 2012*, 2012.
- [8] Bilge Mutlu, Toshiyuki Shiwa and Takayuki Kanda, Hiroshi Ishiguro, and Norihiro Hagita. Footing in human-robot conversations: how robots might shape participant roles using gaze cues. In *Proceedings of the 4th ACM/IEEE International Conference on Human Robot Interaction*, pages 61–68, 2009.
- [9] Kazuhiro Nakadai, Toru Takahashi, Hiroshi G. Okuno, Hirofumi Nakajima, Yuji Hasegawa, and Hiroshi Tsujino. Design and implementation of robot audition system 'HARK' - open source software for listening to three simultaneous speakers. *Advanced Robotics*, 5:739–761, 2010.
- [10] Ralph O. Schmidt. Multiple emitter location and signal parameter estimation. *IEEE Transactions on Antennas and Propagation*, 34:276 – 280, 1986.
- [11] Rainer Stiefelhage, Jie Yang, and Alex Waibel. Modeling focus of attention for meeting indexing based on multiple cues. *IEEE Transactions on Neural Networks*, 13:928–938, 2002.
- [12] David Traum and Jeff Rickel. Embodied agents for multi-party dialogue in immersive virtual worlds. In *Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent Systems*, pages 766 – 773, 2002.