

マイクロホンアレイを用いたサウンドマスキングシステムの検討 —音声入力において指向性を利用したマスキング音レベル低減の試み—

☆小幡将信（上智大院），△秋山あい（上智大），△生田萌人（上智大院），
日岡裕輔（オークランド大），荒井隆行（上智大）

1 はじめに

薬局や病院の待合室，学校の保健室，銀行の窓口などコミュニケーションが不可欠な場面におけるプライベートな内容に対し，プライバシーや機密内容を保護する取り組み，スピーチプライバシーに関する研究が50年以上前から欧米で行われている[1]．こうした取り組みの中でマスキング音を周囲に流すことで，漏れてくる会話内容を聞き取りにくくする技術はサウンドマスキングと呼ばれる．しかし，サウンドマスキングの問題としてマスキング音による第三者への会話の妨げや不快感の増加があげられる[2]．そのため，ターゲット音に対しマスキング音は必要以上の音量を流さないことが好ましい．

本研究ではこの点を踏まえ，想定聴取状況において，第三者へ効率的にマスキング音を流すマスキングシステムの開発を試み，作成したシステムの有効性について考察を行う．

2 問題設定と仮説

2.1 問題設定

診察室で医者と患者が会話をしている間，外の待合室では次の診察を待っている患者がおり，医者との会話を待機者には聞かせたくないという状況を例にする．図1のように，患者を話者A，医者を話者B，マスキング音を流すスピーカーC，D，待機者を聴取位置E，Fと定める．話者Aと話者Bは交互に発話し，マスキング音をスピーカーC，Dから流し聴取位置E，Fにおいてその会話をマスクしたい．

その状況において，スピーカーC，Dから流すマスキング音のレベルをどこまで下げられるかという問題を考える．

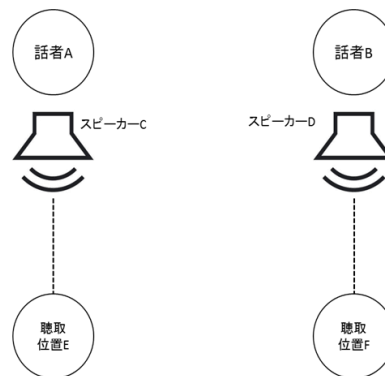


図1 想定聴取状況

2.2 仮説

N. Marrone ら(2008)の研究[3]において，スピーカーの配置によるマスキング効果の違いが提唱された．この研究から「マスキング音の再生位置が変わるとマスキング効果も変わる．特に，話者位置—マスキング音—聴取位置が直線上にあるとき一番マスキング効果が高くなる」ということがわかる．このことから，2.1.で想定した状況において話者Aが発話している間，話者A—スピーカーC—聴取位置Eは直線上にあるため，聴取位置EにおいてスピーカーCのマスキング音が最も影響を与えると考えられる．また，話者A—聴取位置Fの直線上に最も近いスピーカーはスピーカーCであるため，聴取位置FにおいてもスピーカーCのマスキング音が最も影響を与えると考えられる．

同様にして話者Bが発話している間，聴取位置E，FともにスピーカーDのマスキング音が最も影響を与えると考えられる．

よって，話者Aが発話している間は主にスピーカーCからマスキング音が流れ，逆に話者Bが発話している間は主にスピーカーDからマスキング音が流れるという形にすることで，全体としてマスキング音の音量を低下させつつマスキング効果も損なわないモデルができるのではないかと考えた．

“Investigating the use of microphone array for sound masking system: Towards reducing masking sound level using the directivity of input sounds,” Mochinobu Obata, Ai Akiyama, Moeto Ikuta (Sophia Univ.), Yusuke Hioka (Univ. of Auckland), Takayuki Arai (Sophia Univ.)

以上を実現するにあたり、マイクロホンアレイを利用した音声分離を用いた方法を提案する。

3 提案手法 (システムモデル)

仮説に基づき、それぞれの話者において異なるマスクング音を流すシステムを作成するため、マイクロホンアレイによる音声分離と、それを用いた話者の音量変化に伴い音量レベルが変化するマスクング音[4]の作成及び効果的なスピーカー配置による再生を行う。ただし、本研究ではリアルタイムで録音を行うのではなく、オフラインでの録音データを使う。

3.1 音声分離

図2で示すように話者Aと話者Bの音声をそれぞれソース音源 s_1, s_2 、その音声がマイクロホンアレイに届くまでの伝播経路の伝達関数 $a(\theta)$ で表わすとき、 P 個のマイクロホンからなるマイクロホンアレイで録音した音声を録音データ x_p ($p = 1, 2, \dots, P$)に対し任意の角度に指向性に向けたビームフォームのフィルタ係数 $w_p(\theta)$ ($p = 1, 2, \dots, P$)得られた出力を y とする。ただし、本研究ではマイク数 $P = 4$ で進めており、図2においても $P = 4$ とする。

ここでは日岡ら(2013)によって提案された手法[5]を元に、遅延和ビームフォーマとウィナーフィルタを用いて、録音データ x_1, x_2, \dots, x_p からソース音源 s_1, s_2 の音声分離を行った。

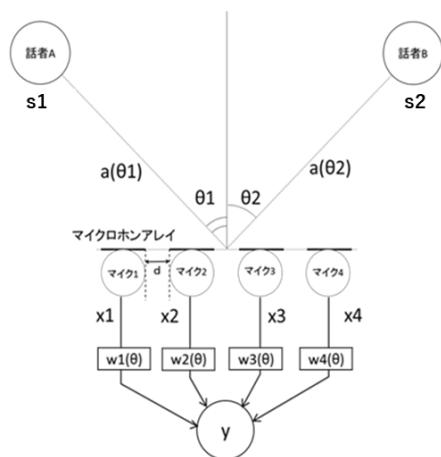


図2 信号とフィルタ

3.2 マスクング音作成

マスクング音を作成するにあたり川のせせらぎ音をオリジナルとなる音として使用する。

これを M_0 とする。また3.1の音源分離で得た、話者A方向から取り出した音声 s_1 を使用したマスクング音、話者B方向から取り出した音声 s_2 を使用したマスクング音、マイク1単体で話者A, Bの会話を収録した信号 x_1 を使用したマスクング音の3種類を作成する。これらをそれぞれ M_A, M_B, M_T とする。これらの作成手順は後の段落に記述する。

先行研究(例えば[6])において音声マスカが用いられているが、その特徴の1つとしてレベルが追従する点があげられる。本研究では、この追従型のマスクング音が作成しやすいと考え採用した。これ以降 M_A を例にして作成法を示していく。

音源信号 $s(\theta_1)$ の最初の2秒間のサンプル数を N とし、各サンプルの二乗和を U とする。 U を次の式(2)に示す。

$$U = \sum_{i=1}^N |s_1(i)|^2 \quad (2)$$

M_A の最初の2秒間にはそのまま M_0 の最初の2秒間をあてはめる。これを式(3)に示す。

$$M_A(i) = M_0(i) \quad (i = 1 \sim N) \quad (3)$$

その後は、サンプルごとに U と過去2秒サンプルの二乗和 q_i の比と忘却定数 α を使用して以下のように求めていく。これを式(4), (5)に示す。

$$q_i = \sum_{j=i-N}^i |s_1(j)|^2 \quad (4)$$

$(i = N + 1 \sim \text{length}(M_0))$

$$M_A(i) = \alpha \times M_A(i-1) + (1 - \alpha) \times \frac{q_i}{U} \times M_0(i) \quad (5)$$

$(i = N + 1 \sim \text{length}(M_0), \text{忘却定数} \alpha = 0.5)$

同様に M_B, M_T も作成する。

4 従来法との比較

今回のマイクロホンアレイを使用した提案手法では以下の二点を利用することで性能の改善を試みている。

- ①音声分離で抽出した音の情報,
- ②目的音との位置関係を考慮した、最適な位置にあるスピーカーの選択

そのため、音声分離を行わない従来の単マイクで集音したデータのマスクング音を二つのスピーカーから同じ音量レベルの音で流した場合を従来法とする。

4.1 想定比較状況

図1で示した状況において、提案法として話者 A, B の会話をマイクロホンアレイにより取り出した音声音源信号 s_1, s_2 を使用して作成したマスキング音 M_A, M_B をそれぞれスピーカーC, D から流すパターンと、従来法として話者 A, B の会話を単マイクで録音した音声 x_1 を使用して作成したマスキング音 M_T をスピーカーC, D 同時に流すパターンを比較する。

4.2 録音実験

まず、図3で示すように4つのマイクロホンアレイ話者 A と話者 B の座標をそれぞれ定めた。

次に、話者 A と話者 B の位置から流したインパルス音を OCTA-CAPTURE を使用して4つのマイクロホンアレイで録音した。また、事前に21歳男性と23歳女性の音声をそれぞれ話者 A, B の発話音声として録音した。これらの音声をたたみ込むことでシミュレーション音声を作成し[7]、話者 A と話者 B の会話をマイクロホンアレイで録音した信号 x_1, x_2, x_3, x_4 として使用した。

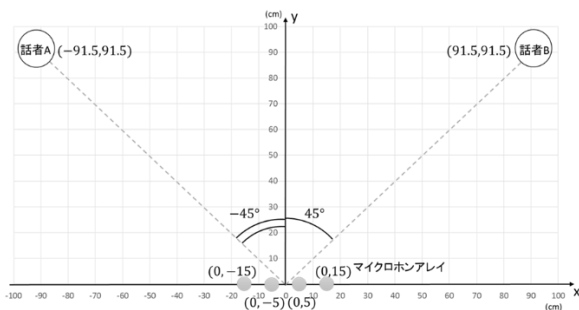


図3 話者とマイクの座標

4.3 結果

4.2の実験から得た録音データ x_1, x_2, x_3, x_4 を使用し、3.1.で示した分離過程を経て得られた音源信号 s_1, s_2 、3.2.で示した作成過程を経て得られたマスキング音 M_A, M_B, M_T 、比較のための単マイク録音信号 x_1 の波形を図5に示す。ただし、音声波形の横軸の単位は $\times 10^5$ サンプルである(サンプリング周波数は16 kHz)。

それぞれの波形の形から、音源信号 s_1, s_2 や単マイクで録音した音声信号 x_1 の振幅が大きくなったとき、マスキング音 M_A, M_B, M_T の振

幅も追従して大きくなっていることがわかる。

また、音声分析ソフト Pratt を使用し、それぞれのマスキング音のインテンシティの平均レベルを測ると、 M_A が約 58.0 dB, M_B が約 57.0 dB であるのに対し、 M_T は約 57.3 dB となった。近い値になったが、想定状況において従来法では両方のスピーカーから M_T が常時流れるのに対し、提案法ではスピーカーが交互に音が流れる形となっているため総合的にみると従来法のレベルが大きいと考えられる。

さらに、図5に s_1 の音声波形と M_A の音声波形を重ね、話者 A が発話する一部分を拡大した音声波形を示す。

この図のより話者 A が発話してからマスキング音が十分なレベルになるまで約 1500 サンプル、つまり約 0.1 秒の遅延が発生していることがわかる。

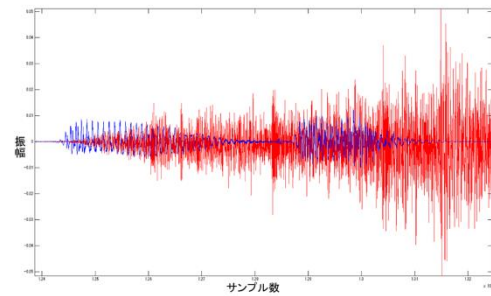


図4 s_1 と M_A を重ねた音声波形と一部拡大したもの

5 考察

4.3.で得られた結果から、4.1.で想定した状況においてマイクロホンアレイを使用することで話者 A が発話している間はスピーカーCから、話者 B が発話している間はスピーカーDから発話者の音声に合わせた音量レベルのマスキング音が流れる状況ができており、これは仮説で考えたモデルが実現できていると考えられる。

話者が発話を開始したタイミングにおいて遅延の影響でマスキング音のレベルが十分になっていないという問題については、話者の発話を感知した瞬間にある程度のレベルのマスキング音を流し、そこからレベル変化の調整をしていくというシステムによって改善できると予想される。

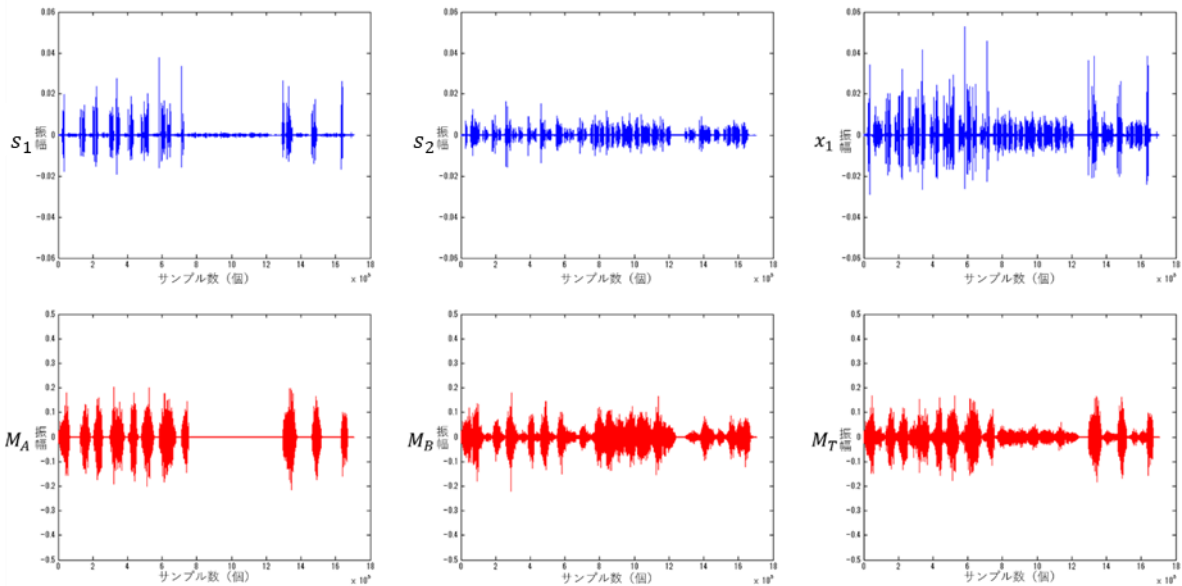


図 5 音源信号とマスクング音の音声波形

6 おわりに

本研究では、マイクロホンアレイを用いて音声分離を行い、マスクング音を作成することでマスクング効果を損なわずに音量低下も見込めるようなシステムについて提案し、その可能性を探った。

今後は、聴取者を対象にしてマスクング効果を評価するために明瞭度や Annoyance を測る聴取実験を行う必要がある。また、今回考えた仮説は定量的な議論が出来ていないため、それを考慮したシステムモデルも今後検討していく。

謝辞

荒井研究室の Justine Hui さんには、実験の指導や論文の執筆等アドバイスしていただき感謝申し上げます。

参考文献

- [1] W. J. Cavanaugh, W. R. Farrell, P. W. Hirtle and B. G. Watters, "Speech privacy in buildings," *J. Acoust. Soc. Am.*, 34, pp. 475–492, 1962.
- [2] 王循, 藤田佑一郎, 木庭洋介, 石川諭, 雉本信哉, "各種マスキングによる音声マスクング効果及び心理的影響の比較," 第 25 回環境工学総合シンポジウム論文集, pp. 14–17, 2015.

- [3] N. Marrone, C. R. Mason, G. Kidd, Jr, "Tuning in the spatial dimension: Evidence from a masked speech identification task," *J. Acoust. Soc. Am.*, 124, pp.1146–1158, 2008.
- [4] 李孝珍, 上野佳奈子, 坂本慎, 藤原舞, 清水寧, "レベル適応型マスクングシステムの有効性に関する検討," 日本音響学会研究発表会講演論文集, pp. 1077–1080, 2010.
- [5] Y. Hioka, K. Furuya, K. Kobayashi, K. Niwa and Y. Haneda, "Underdetermined Sound Source Separation Using Power Spectrum Density Estimated by Combination of Directivity Gain," in *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 6, pp. 1240–1250, 2013.
- [6] T. Arai, "Masking speech with its time-reversed signal," *Acoust. Sci. & Tech.*, 31(2), pp. 188–190, 2010.
- [7] 橋秀樹, 日高新人, "実物及び模型ホールのインパルス応答の測定," *日本音響学会誌*, 48 巻, 4 号, pp.244–249, 1992.