

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5178370号  
(P5178370)

(45) 発行日 平成25年4月10日 (2013. 4. 10)

(24) 登録日 平成25年1月18日 (2013.1.18)

(51) Int. Cl.	F I
G 1 O L 21/0308 (2013.01)	G 1 O L 21/02 2 O 3 Z
G 1 O L 21/028 (2013.01)	G 1 O L 21/02 2 O 1 D
G 1 O L 21/0272 (2013.01)	G 1 O L 21/02 2 O 2 B

請求項の数 2 (全 12 頁)

(21) 出願番号	特願2008-191382 (P2008-191382)	(73) 特許権者	000005326
(22) 出願日	平成20年7月24日 (2008. 7. 24)		本田技研工業株式会社
(65) 公開番号	特開2009-42754 (P2009-42754A)		東京都港区南青山二丁目1番1号
(43) 公開日	平成21年2月26日 (2009. 2. 26)	(74) 代理人	110000800
審査請求日	平成22年11月26日 (2010. 11. 26)		特許業務法人創成国際特許事務所
(31) 優先権主張番号	60/954889	(72) 発明者	武田 龍
(32) 優先日	平成19年8月9日 (2007. 8. 9)		埼玉県和光市本町8-1 株式会社ホンダ
(33) 優先権主張国	米国 (US)		・リサーチ・インスティテュート・ジャパ ン内
		(72) 発明者	中臺 一博
			埼玉県和光市本町8-1 株式会社ホンダ
			・リサーチ・インスティテュート・ジャパ ン内

最終頁に続く

(54) 【発明の名称】 音源分離システム

(57) 【特許請求の範囲】

【請求項1】

スピーカから環境に対して音響として出力される既知信号を保存する既知信号記憶手段と、

マイクロホンと、

前記マイクロホンからの出力信号を周波数変換することにより現在のフレームの観測信号を生成する第1処理部と、

周波数領域における前記既知信号の周波数成分と前記スピーカから前記マイクロホンまでの周波数領域の音響の伝達関数との畳み込みとして現在フレームの元信号が表現されている第1モデルと、前記元信号および未知信号を包含するように前記観測信号が表現されている第2モデルとにしたがって、前記第1処理部により生成された現在フレームの前記観測信号を、適応フィルタを經由した前記元信号と前記未知信号とに独立成分分析法により分離することにより前記未知信号を抽出する第2処理部とを備えていることを特徴とする音源分離システム。

【請求項2】

請求項1記載の音源分離システムにおいて、

前記適応フィルタ  $h^{\wedge}$  の更新式は次の (a) 式とし、

$$h^{\wedge}(\quad, f + 1) = h^{\wedge}(\quad, f) - \mu_1 (E(\quad, f)) X^*(\quad, f) \cdots (a)$$

とし、

ただし、前記 (a) 式において、

は周波数を意味し、  
 $f + 1$  ,  $f$  は更新後および更新前のフレームを意味し、  
 $\mu_1$  は所定の係数であり、  
 $E(\cdot, f)$  は前記未知信号であり、  
関数 は確率変数  $e$  の密度関数  $p_x(x)$  により次の (b) 式により定義されるものであり、

$$p_x(x) = - (d / dx) \log p_x(x) \cdot \dots \cdot (b)$$

$X^*(\cdot, f)$  は前記元信号  $X(\cdot, f)$  の複素共役であることを特徴とする音源分離システム。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明は、音源分離システムに関する。

【背景技術】

【0002】

ユーザとロボットとの自然な対話を実現する上で、ロボットの発話中にユーザの発話（いわゆるバージイン）を許容することは不可欠である。ロボットにマイクロホンが搭載されている場合、ロボット自身の発話がマイクロホンに入り込むので、バージインは相手の発話を認識する上で大きな障害となる。

【0003】

20

そこで、図4に示されている構成の適応フィルタが用いられている。自発話の除去は、スピーカSからマイクロホンMへの伝達系  $h$  を近似するフィルタ  $h^{\wedge}$  の推定問題として取り扱われる。マイクロホンMから入力された観測信号  $y(k)$  から推定信号  $y^{\wedge}(k)$  が差し引かれることにより相手発話を取り出される。

【0004】

適応フィルタの1つとしてNLS (Normalized Least Mean Squares) 法が提案されている。NLS法によれば、時間領域において、線形かつ時間不変な伝達系を経て観測される信号  $y(k)$  が、元信号ベクトル  $x(k) = {}^t(x(k), x(k-1), \dots, x(k-N+1))$  (「N」はフィルタ長を表わす。「t」は転置を意味する。) と、伝達系のインパルス応答  $h = {}^t(h_1, h_2, \dots, h_N)$  との畳み込みを用いて関係式(1)により表現される。

30

【0005】

$$y(k) = {}^t x(k) h \dots (1)$$

【0006】

推定フィルタ  $h^{\wedge} = {}^t(h_1^{\wedge}, h_2^{\wedge}, \dots, h_N^{\wedge})$  は、関係式(2)により表わされる観測信号と推定信号との誤差  $e(k)$  の二乗平均を最小化することにより得られる。推定フィルタ  $h^{\wedge}$  を求めるためのオンラインアルゴリズムは、正則化のための小さな正数値を用いて関係式(3)により表現される。なお、関係式(3)において  $\|x(k)\|^2 + \epsilon$  により学習係数が正規化されない場合がNLS法である。

【0007】

$$e(k) = y(k) - {}^t x(k) h^{\wedge} \dots (2)$$

40

【0008】

$$h^{\wedge}(k) = h^{\wedge}(k-1) + \mu_{NLS} x(k) e(k) / (\|x(k)\|^2 + \epsilon) \dots (3)$$

【0009】

また、ICA (独立成分分析) 法が提案されている。ICA法はノイズを仮定して設計されているため、自発話区間の検出が不要であり、かつ、ノイズが存在しても分離可能であり、その結果としてバージイン問題の解決に適している。たとえば、時間領域ICA法が提案されている(非特許文献1参照)。音源の混合過程はノイズ  $n(k)$  および  $N+1$  次の行列  $A$  を用いて関係式(4)により表わされる。

【0010】

$${}^t(y(k), {}^t x(k)) = A {}^t(n(k), {}^t x(k)),$$

50

$$A_{ii} = 1 (i = 1 \dots N + 1), A_{1j} = h_{j-1} (j = 2 \dots N + 1), \\ A_{ik} = 0 (k \neq i) \dots (4)$$

【 0 0 1 1 】

ICAによれば関係式(5)における分離行列Wが推定される。

【 0 0 1 2 】

$${}^t(e(k), {}^t x(k)) = W^{-1}({}^t y(k), {}^t x(k)), \\ W_{11} = a, W_{ii} = 1 (i = 2 \dots N + 1), \\ W_{1j} = h_j (j = 2 \dots N + 1), W_{ik} = 0 (k \neq i) \dots (5)$$

【 0 0 1 3 】

分離行列Wの第1行第1列成分 $W_{11} = a = 1$ の場合が従来の適応フィルタのモデルであり、ICA法と最も異なっている点である。自然勾配法にしたがってKL情報量が最小化されることにより、オンラインアルゴリズムを表わす関係式(6)および(7)にしたがって最適な分離フィルタが求められる。

10

【 0 0 1 4 】

$$h^{(k+1)} = h^{(k)} + \mu_1 [ \{ 1 - (e(k))e(k) \} h^{(k)} - (e(k))x(k) ] \dots (6)$$

【 0 0 1 5 】

$$a^{(k+1)} = a^{(k)} + \mu_2 [ 1 - (e(k))e(k) ] a^{(k)} \dots (7)$$

【 0 0 1 6 】

関数  $\ln p_x(x)$  は確率変数  $e$  の密度関数  $p_x(x)$  により関係式(8)により定義される。

20

【 0 0 1 7 】

$$\ln p_x(x) = - \int (d/dx) \ln p_x(x) \dots (8)$$

【 0 0 1 8 】

さらに、周波数領域ICA法が提案されている(非特許文献2参照)。一般的に周波数領域では畳み込み混合が瞬時の混合とみなせるため、時間領域ICA法よりも収束性に優れている。この手法によれば、窓長Tおよびシフト長Uによる短時間フーリエ解析が実行されることにより、時間周波数領域での信号が得られる。元信号 $x(t)$ および観測信号 $y(t)$ のそれぞれはフレームfおよび周波数  $f$  を変数とする関数 $X(\cdot, f)$ および $Y(\cdot, f)$ のそれぞれにより表現される。観測信号ベクトル $Y(\cdot, f) = {}^t(Y(\cdot, f), X(\cdot, f))$ の分離過程は、推定された元信号ベクトル $Y^{\wedge}(\cdot, f) = {}^t(E(\cdot, f), X(\cdot, f))$ を用いて関係式(9)により表現される。

30

【 0 0 1 9 】

$$Y^{\wedge}(\cdot, f) = W(\cdot) Y(\cdot, f), W_{21}(\cdot) = 0, W_{22}(\cdot) = 1 \dots (9)$$

【 0 0 2 0 】

分離行列の学習は周波数ごとに独立に行われる。学習は非ホロノミック拘束適用によるKL情報量最小化に基づく関係式(10)により表わされる反復学習則にしたがって行われる(非特許文献3参照)。

【 0 0 2 1 】

$$W^{(j+1)}(\cdot) = W^{(j)}(\cdot) - \{ \text{off-diag} \langle Y^{\wedge} Y^{\wedge H} \rangle \} W^{(j)}(\cdot) \dots (10)$$

【 0 0 2 2 】

$\alpha$  は学習係数であり、 $(j)$  は更新回数であり、 $\langle \cdot \rangle$  は平均値であり、 $\text{off-diag} X$  は行列Xの対角要素を0に置換する演算を表わし、非線形関数  $\tanh(y)$  は関係式(11)により定義されている。

40

【 0 0 2 3 】

$$\tanh(y_i) = \tanh(|y_i|) \exp(i \arg(y_i)) \dots (11)$$

【 0 0 2 4 】

また、既知音源から既知音源への伝達特性は定数で表わされるため、分離行列Wの第1行成分のみが更新される。

【非特許文献1】J.Yang et al., A New Adaptive Filter Algorithm for System Identification Using Independent Component Analysis, Proc. ICASSP2007, pp.1341-1344, 20

50

07

【非特許文献2】S. Myabe et al., Double-Talk Free Spoken Dialogue Interface Combining Sound Field Control with Semi-Blind Source Separation, Proc. ICASSP2006, p.809-812, 2006

【非特許文献3】Sawada et al., Polar Coordinate based Nonlinear Function for Frequency-Domain Blind Source Separation, IEICE Trans., Fundamentals, 3, E-86A, pp. 505-510, 2003

【発明の開示】

【発明が解決しようとする課題】

【0025】

しかし、従来の周波数領域ICA法には次のような問題点があった。第1の問題は残響に対応するために窓長Tを長く取る必要があり、その分だけ演算処理遅延および音源分離性能の低下を招くという点である。第2の問題点は窓長Tを環境に応じて変更する必要があり、他の雑音抑圧手法などの接続が煩雑になる点である。

【0026】

そこで、本発明は、残響または反響の影響を軽減することにより音源分離精度の向上を図ることができるシステムを提供することを課題とする。

【課題を解決するための手段】

【0027】

第1発明の音源分離システムは、スピーカから環境に対して音響として出力される既知信号を保存する既知信号記憶手段と、マイクロホンと、前記マイクロホンからの出力信号を周波数変換することにより現在のフレームの観測信号を生成する第1処理部と、周波数領域における前記既知信号の周波数成分と前記スピーカから前記マイクロホンまでの周波数領域の音響の伝達関数との畳み込みとして現在フレームの元信号が表現されている第1モデルと、前記元信号および未知信号を包含するように前記観測信号が表現されている第2モデルとにしたがって、前記第1処理部により生成された現在フレームの前記観測信号を、適応フィルタを経由した前記元信号と前記未知信号とに独立成分分析法により分離することにより前記未知信号を抽出する第2処理部とを備えていることを特徴とする。

【0028】

第1発明の音源分離システムによれば、第1モデルおよび第2モデルにしたがって観測信号から未知信号が抽出される。特に、第1モデルによれば、現在フレームの元信号が現在および過去フレームの既知信号の合成信号として表現されている。このため、窓長を変更させることなく、既知信号の残響または反響が観測信号に及ぼす影響を軽減しながら未知信号が抽出されうる。したがって、残響の影響を軽減するための演算処理負荷を軽減しながら、未知信号に基づく音源分離精度の向上を図ることができる。

【0029】

第1発明の音源分離システムによれば、現在フレームの元信号が周波数領域における既知信号の周波数成分およびその伝達関数の畳み込みにより表現されている。このため、窓長を変更させることなく、元信号の残響または反響が観測信号に及ぼす影響を軽減しながら未知信号が抽出されうる。したがって、残響の影響を軽減するための演算処理負荷を軽減しながら、未知信号に基づく音源分離精度の向上を図ることができる。

【0030】

第1発明の音源分離システムによれば、第2モデルにおいて適応的に分離フィルタが設定されるので、窓長を変更させることなく、元信号の残響または反響が観測信号に及ぼす影響を軽減しながら未知信号が抽出されうる。したがって、残響の影響を軽減するための演算処理負荷を軽減しながら、未知信号に基づく音源分離精度の向上を図ることができる。

【0031】

第2発明の音源分離システムは、第1発明の音源分離システムにおいて、前記適応フィルタ $h^{\wedge}$ の更新式は次の(a)式とし、

10

20

30

40

50

$$h^{\wedge}(\quad, f + 1) = h^{\wedge}(\quad, f) - \mu_1 (E(\quad, f)) X^{\wedge}(\quad, f) \cdot \cdot \cdot (a)$$

とし、ただし、前記 (a) 式において、 $\quad$  は周波数を意味し、 $f + 1$ 、 $f$  は更新後および更新前のフレームを意味し、 $\mu_1$  は所定の係数であり、 $E(\quad, f)$  は前記未知信号であり、関数  $\quad$  は確率変数  $e$  の密度関数  $p_x(x)$  により次の (b) 式により定義されるものであり、

$$(\quad) = - (d / dx) \log p_x(x) \cdot \cdot \cdot (b)$$

$X^{\wedge}(\quad, f)$  は前記元信号  $X(\quad, f)$  の複素共役であることを特徴とする。

【発明を実施するための最良の形態】

【0033】

本発明の音源分離システムの実施形態について図面を用いて説明する。

10

【0034】

図1に示されている音源分離システムはマイクロホンMと、スピーカSと、電子制御ユニット(CPU, ROM, RAM/O回路、A/D変換回路等の電子回路などにより構成されている。)10とにより構成されている。電子制御ユニット10は第1処理部11と、第2処理部12と、第1モデル格納部101と、第2モデル格納部102と、自発話格納部104とを備えている。各処理部はたとえば演算処理回路、または、メモリと、メモリからプログラムを読み出してそのプログラムにしたがって担当する演算処理を実行する演算処理装置(CPU)とにより構成されている。

【0035】

第1処理部11はマイクロホンMからの出力信号を周波数変換することにより現在のフレーム  $f$  の観測信号(周波数成分)  $Y(\quad, f)$  を生成する。第2処理部12は第1処理部11により生成された現在フレームの観測信号  $Y(\quad, f)$  に基づき、第1モデル格納部101に格納されている第1モデルと、第2モデル格納部102に格納されている第2モデルとにしたがって未知信号  $E(\quad, f)$  を抽出する。電子制御ユニット10は自発話格納部(既知信号記憶手段)104に格納されている既知信号をスピーカSから音声または音響として出力させる。

20

【0036】

マイクロホンMはたとえば図2に示されているように電子制御ユニット10が搭載されているロボットRの頭部P1に配置されている。なお、音源分離システムはロボットRのほか、車両(四輪自動車)、複数の音源が存在する環境に接する任意の機械や装置に搭載されうる。また、マイクロホンMの数および配置は任意に変更されうる。ロボットRは脚式移動ロボットであり、人間と同様に基体P0と、基体P0の上方に配置された頭部P1と、基体P0の上部に上部両側から延設された左右の腕体P2と、左右の腕体P2のそれぞれの先端に連結されている手部P3と、基体P0の下部から下方に延設された左右の脚体P4と、左右の脚体P4のそれぞれに連結されている足部P5とを備えている。基体P0はヨー軸回りに相対的に回動しうるように上下に連結された上部および下部により構成されている。頭部P1は基体P0に対してヨー軸回りに回動する等、動くことができる。腕体P2は肩関節機構、肘関節機構および手根関節機構のそれぞれにおいて1~3軸回りの回動自由度を有している、手部P3は、手掌部から延設され、人間の手の親指、人差指、中指、薬指および小指のそれぞれに相当する5つの指機構を備え、物体の把持動作等が可能に構成されている。脚体P4は股関節機構、膝関節機構および足関節機構のそれぞれにおいて1~3軸回りの回動自由度を有している。ロボットRは音源分離システムによる音源分離結果に基づき、左右の脚体P4を動かして移動する等、適当な動作をすることができる。

30

40

【0037】

前記構成の音源分離システムの機能について説明する。まず第1処理部11によりマイクロホンMからの出力信号が取得される(図3/S002)。また、第1処理部11によりこの出力信号がA/D変換された上で周波数変換されることにより、フレーム  $f$  の観測信号  $Y(\quad, f)$  が生成される(図3/S004)。

【0038】

50

続いて第2処理部12により、第1モデルおよび第2モデルにしたがって、第1処理部11により生成された観測信号 $Y(n, f)$ から元信号 $X(n, f)$ が分離されることにより、未知信号 $E(n, f)$ が抽出される(図3/S006)。

【0039】

第1モデルによれば、現在および過去の所定数 $M$ のフレームにわたる元信号を包含するように現在フレーム $f$ の元信号 $X(n, f)$ が表現されている。第1モデルによれば、次フレームに入り込んだ反響音が時間周波数領域における畳み込みにより表現されている。具体的には、あるフレーム $f$ の周波数成分が $M$ フレームにわたって観測信号の周波数成分に影響を及ぼすという仮定のもと、元信号 $X(n, f)$ が、遅延した既知信号(具体的には元信号の遅延 $m$ の周波数成分) $S(n, f - m + 1)$ およびその伝達関数 $A(n, m)$ の畳み込みとして関係式(12)により表現されている。

【0040】

$$X(n, f) = \sum_{m=1}^M A(n, m) S(n, f - m + 1) \dots (12)$$

【0041】

図5には当該畳み込みの模式図が示されている。畳み込まれた未知信号 $E(n, f)$ と、通常の伝達過程を経た既知音(自発話信号) $S(n, f)$ との混合が観測音 $Y(n, f)$ であるとみなされる。これは、一様DTF(Discrete Fourier Transform)フィルタバンクによる一種のマルチレート処理に相当する。

【0042】

第2モデルによれば、適応フィルタ(分離フィルタ) $h^{\wedge}$ を経由した元信号 $X(n, f)$ と観測信号 $Y(n, f)$ とを包含するように未知信号 $E(n, f)$ が表現されている。具体的には、第2モデルによる分離過程は、元信号ベクトル $X$ 、未知信号 $E$ 、観測音スペクトル $Y$ および分離フィルタ $h^{\wedge}$ および $c$ に基づき、関係式(13)~(15)にしたがってベクトル表現されている。

【0043】

$$\begin{aligned} {}^t(E(n, f), {}^tX(n, f)) &= C {}^t(Y(n, f), {}^tX(n, f)), \\ C_{11} &= c(n), C_{ii} = 1 (i = 2 \dots M + 1), \\ C_{1j} &= h_{j-1}^{\wedge} (j = 2 \dots M + 1), C_{ki} = 0 (k \neq i) \dots (13) \end{aligned}$$

【0044】

$$X(n, f) = {}^t(X(n, f), X(n, f - 1), \dots, X(n, f - M + 1)) \dots (14)$$

【0045】

$$h^{\wedge}(n) = (h_1^{\wedge}(n), h_2^{\wedge}(n) \dots h_M^{\wedge}(n)) \dots (15)$$

【0046】

この表現は複素数が用いられるほかは時間領域ICA法と同一であるが、収束性の観点から周波数領域ICA法においてよく利用されている関係式(11)が用いられた。これによりフィルタ $h^{\wedge}$ の更新は関係式(16)により表現される。

【0047】

$$h^{\wedge}(f + 1) = h^{\wedge}(f) - \mu_1 (E(f)) X^{\wedge}(f) \dots (16)$$

【0048】

$X^{\wedge}(f)$ は $X(f)$ の複素共役を表わしている。なお、周波数インデックスは省略されている。

【0049】

分離フィルタ $c$ に関する更新がないため、分離フィルタ $c$ は分離行列の初期値 $c_0$ のままである。初期値 $c_0$ は誤差 $E$ の対数密度関数の導関数 $\psi(x)$ に対して適切に定められるスケール係数である。関係式(16)から明らかなようにフィルタ更新時に誤差(未知信号) $E$ が適切にスケールされていれば、その学習は阻害されない。このため、スケール係数 $a$ がなんらかの方法にしたがって求められ、これを用いて関数 $\psi(aE)$ が適用されれば、分離行列の初期値 $c_0$ が1であっても差し支えない。スケール係数の学習則は時間領域ICA法と同様の関係式(7)が用いられればよい。これは、関係式(7)によれば実質的に $e$ を正規化するスケール係数が求められているからである。時

10

20

30

40

50

間領域 I C A 法における  $e$  は  $a E$  に相当する。

【 0 0 5 0 】

以上から第 2 モデルによる学習則は関係式 ( 1 7 ) ~ ( 1 9 ) により表現される。

【 0 0 5 1 】

$$E(f) = Y(f) - {}^t X(f) h^{\wedge}(f) \dots (17)$$

【 0 0 5 2 】

$$h^{\wedge}(f+1) = h^{\wedge}(f) + \mu_1 (a(f) E(f)) X^*(f) \dots (18)$$

【 0 0 5 3 】

$$a(f+1) = a(f) + \mu_2 [1 - (a(k) E(k)) a^*(f) E^*(f)] a(f) \dots (19)$$

10

【 0 0 5 4 】

非線形関数  $\sigma(x)$  が  $\tanh(|x|) \exp(i \phi(x))$  等、 $r(|x|, \phi(x)) \exp(i \phi(x))$  という形式を満たしていれば  $a$  は実数となる。

【 0 0 5 5 】

前記機能を発揮する音源分離システムによれば、第 1 モデルおよび第 2 モデルにしたがって、観測信号  $Y(n, f)$  から未知信号  $E(n, f)$  が抽出される ( 図 3 / S 0 0 2 ~ S 0 0 6 参照 )。第 1 モデルによれば、現在フレーム  $f$  の観測信号  $Y(n, f)$  が現在および過去の所定数  $M$  のフレームにわたる元信号  $X(n, f - m + 1)$  ( $m = 1 \sim M$ ) の合成信号として表現されている ( 関係式 ( 1 2 ) 参照 )。また、第 2 モデルにおいて適応的に分離フィルタ  $h^{\wedge}$  が設定される ( 関係式 ( 1 6 ) ~ ( 1 9 ) 参照 )。このため、窓長を変更させることなく、元信号  $X(n, f)$  の残響または反響が観測信号  $Y(n, f)$  に及ぼす影響を軽減しながら未知信号  $E(n, f)$  が抽出されうる。したがって、既知信号  $S(n, f)$  の残響の影響を軽減するための演算処理負荷を軽減しながら、未知信号  $E(n, f)$  に基づく音源分離精度の向上を図ることができる。

20

【 0 0 5 6 】

ここで、関係式 ( 3 ) および ( 1 8 ) を比較する。適用領域を除けば、本願発明の拡張周波数領域 I C A 法はスケーリング係数  $a$  および関数  $\sigma$  により L M S ( N L M S ) 法における推定フィルタと相違している。簡単のため、定義域が時間領域 ( 実数 ) であり、ノイズ ( 未知信号 ) が標準正規分布にしたがうと仮定すると、関数  $\sigma$  は関係式 ( 2 0 ) により表わされる。

30

【 0 0 5 7 】

$$\sigma(x) = -(d/dx) \log(\exp(-x^2/2)) / (2)^{1/2} = x \dots (20)$$

【 0 0 5 8 】

これは、関係式 ( 1 8 ) 右辺第 2 項に含まれる  $(a E(t)) X(t)$  が  $a E(t) X(t)$  と表現されることを意味するので、関係式 ( 1 8 ) は関係式 ( 3 ) と等価になる。これは、関係式 ( 3 ) において学習係数が適切に定められれば L M S 法でも Double-Talk 状態においてフィルタ更新が可能であることを意味する。換言すると、ノイズがガウス分布にしたがっており、かつ、学習係数がノイズのパワーに応じて適切に設定されている場合、L M S 法は I C A 法と等価な動作をする。

【 0 0 5 9 】

40

図 6 に L M S 法および I C A 法のそれぞれによる分離例が示されている。観測音は前半では自発話のみである一方、後半では自発話と相手発話とが混じっている。L M S 法によればノイズがない区間では拘束に収束しているが、ノイズがある Double-Talk 状態では不安定な動作を示している。これに対して、I C A 法によれば収束は遅いもののノイズがある区間でも安定である。

【 0 0 6 0 】

続いて、A . 時間領域 N L M S 法、B . 時間領域 I C A 法、C . 周波数領域 I C A 法および D . 本願発明の手法のそれぞれの連続音源分離性能の実験結果について説明する。

【 0 0 6 1 】

実験に際して図 7 に示されているように 4 . 2 m x 7 m の広さの部屋 ( 残響時間 ( R T

50

60) が約 0.3 秒) において、サンプリングレート 16 kHz でインパルス応答が録音された。自発話に対応するスピーカ S はマイク付近に設置され、マイク M に対するスピーカ S が向く方向を正面方向とした。相手発話に対応するスピーカはマイクに向けて設置された。マイク M とスピーカとの距離は 1.5 m とされた。録音されたインパルス応答を積み込んだ ASJ - JNAS の評価用データセット 200 文 (男女各 100 文) が評価用データとして用いられた。この 200 文を相手発話とし、自発話にはその中の一文 (約 7 秒) を用いた。混合されたデータは、相手発話および自発話の始まりは揃っているが終わりは揃っていない。

#### 【0062】

音源分離エンジンとして Julius が使用された (1216856767017\_0 参照)。クリーン音声 200 話者 (男性 100 人、女性 100 人) 分の ASJ - JNAS 新聞記事読み上げ、および、音素バランス文計 150 文で学習したトライフォン (3 状態 8 混合の HMM) が音響モデルとして使用された。MFCC (12 + 12 + Pow) 25 次元が音源分離特徴量として用いられた。認識に用いられた音声は学習データに含まれていない。

#### 【0063】

実験条件を一致させるため、時間領域におけるフィルタ長が約 0.128 秒に設定された。これにより、手法 A および手法 B のフィルタ長は 2048 (約 0.128 秒) となる。手法 D では窓長 T が 1024 (0.064 秒) に設定され、シフト長 U が 128 (約 0.008 秒) に設定され、かつ、遅延フレーム数 M が 8 に設定されることにより手法 A および手法 B と条件を一致させた。手法 C では窓長 T が 2048 (0.128 秒) に設定され、シフト長 U が手法 D と同様に 128 (0.008 秒) に設定された。フィルタの初期値はすべて 0 に設定され、オンライン処理で分離が実行された。

#### 【0064】

学習係数の値としては試行錯誤により認識率が最高になる値が選択された。学習係数は収束性および分離性能を左右する因子であるが、最適値から大きく外れていない限り性能を著しく変化させることはない。

#### 【0065】

図 8 に認識結果である単語認識率が示されている。「観測音」は適応フィルタがない状態、すなわち、なんら処理が施されない状態での認識結果を表わしている。「単独発話」は自発話の混合がない状態、すなわち、ノイズがない状態での認識結果を表わしている。クリーン音声の一般的な認識率は約 90% であることから、図 8 から明らかなように部屋環境の影響によって約 20% も認識率が低下している。手法 A では観測音と比較して認識率が 0.87% だけ低下している。これは、自発話と相手発話とが混在する Double-Talk 状態では手法 A の動作が不安定になることを反映しているためであると推察される。手法 B では観測音と比較して認識率が 4.21% だけ上昇し、手法 C では観測音と比較して認識率が 7.55% だけ上昇している。これは、時間領域で処理が実行される手法 B よりも、周波数領域で処理が実行される結果として周波数ごとの特性が反映される手法 C のほうがよい結果が得られることを表わしている。手法 D では観測音と比較して認識率が 9.61% だけ上昇しており、従来手法 A ~ C よりも有効な音源分離手法であることが確認された。

#### 【図面の簡単な説明】

#### 【0066】

【図 1】本発明の音源分離システムの構成説明図

【図 2】本発明の音源分離システムのロボットへの搭載例示図

【図 3】本発明の音源分離システムの機能を示すフローチャート

【図 4】適応フィルタの構成に関する説明図

【図 5】時間周波数領域における畳み込みに関する説明図

【図 6】LMS 法および ICA 法による相手発話の分離結果に関する説明図

【図 7】実験状況に関する説明図

【図 8】音源分離結果としての各手法による単語正解率の比較説明図

10

20

30

40

50

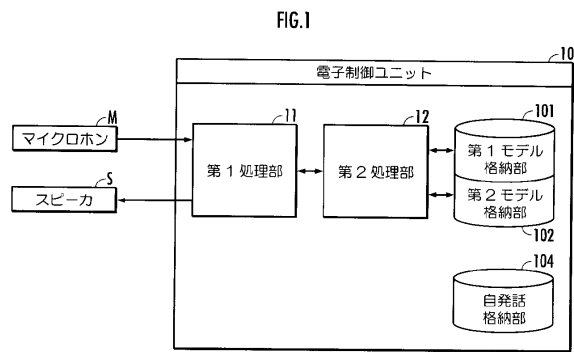


【符号の説明】

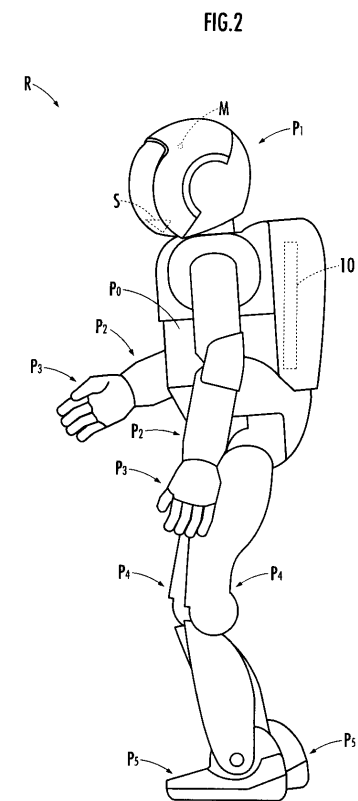
【0067】

10 電子制御ユニット、11 第1処理部、12 第2処理部、S スピーカ、M マイクロホン

【図1】



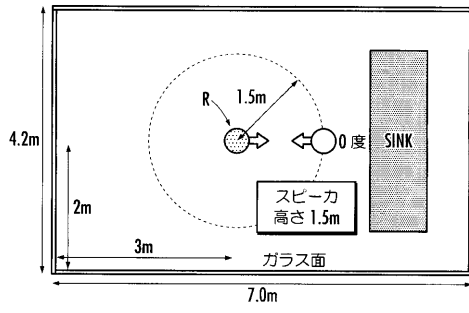
【図2】





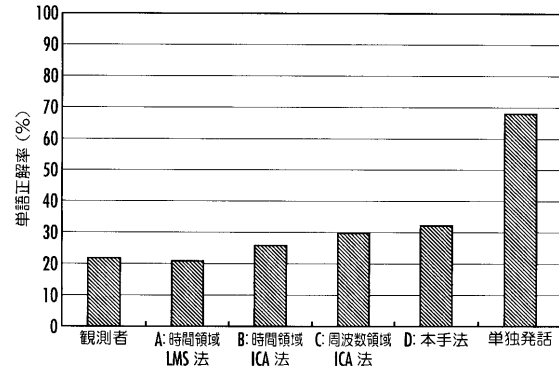
【 図 7 】

FIG.7



【 図 8 】

FIG.8



---

フロントページの続き

(72)発明者 辻野 広司

埼玉県和光市本町8-1 株式会社ホンダ・リサーチ・インスティテュート・ジャパン内

(72)発明者 奥乃 博

京都府京都市中京区東洞院通三条下る三文字町205番地の3-1102

審査官 菊池 智紀

(56)参考文献 特開2006-163231(JP, A)

大田健紘 他, "既知雑音除去法の有効性に関する検討", 電子情報通信学会技術研究報告, 2006年 5月19日, Vol.106, No.78, p.7-12

武田龍 他, "ICAとMFTに基づく音声認識におけるSoft Maskを用いた性能評価", 情報処理学会第69回全国大会予稿集(2), 2007年 3月 6日, p.2-585 2-586

Shigeki Miyabe et al., "Double-talk Free Spoken Dialogue Interface Combining Sound Field Control with Semi-blind Source Separation", Proc. of IEEE ICASSP2006, 2006年 5月14日, Vol.1, p.1-809 1-812

(58)調査した分野(Int.Cl., DB名)

G10L 15/20, 21/02

JSTPlus(JDreamII)