

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5572445号  
(P5572445)

(45) 発行日 平成26年8月13日(2014.8.13)

(24) 登録日 平成26年7月4日(2014.7.4)

(51) Int.Cl. F I  
**G 1 0 L 21/0208 (2013.01)** G 1 0 L 21/0208 1 0 0 B  
**G 1 0 L 21/0272 (2013.01)** G 1 0 L 21/0272 1 0 0 B

請求項の数 8 (全 22 頁)

(21) 出願番号	特願2010-105369 (P2010-105369)	(73) 特許権者	000005326
(22) 出願日	平成22年4月30日 (2010.4.30)		本田技研工業株式会社
(65) 公開番号	特開2011-232691 (P2011-232691A)		東京都港区南青山二丁目1番1号
(43) 公開日	平成23年11月17日 (2011.11.17)	(74) 代理人	100064908
審査請求日	平成24年11月27日 (2012.11.27)		弁理士 志賀 正武
		(74) 代理人	100108578
			弁理士 高橋 詔男
		(74) 代理人	100146835
			弁理士 佐伯 義文
		(74) 代理人	100094400
			弁理士 鈴木 三義
		(74) 代理人	100107836
			弁理士 西 和哉
		(74) 代理人	100108453
			弁理士 村山 靖彦

最終頁に続く

(54) 【発明の名称】 残響抑圧装置、及び残響抑圧方法

(57) 【特許請求の範囲】

【請求項1】

音声信号を取得する音声取得部と、  
 音声信号を出力する音声出力部と、

前記音声出力部から既知の音声信号を出力させたときに、前記音声取得部によって取得された残響音を含む前記既知の音声信号と残響音を含まない当該既知の音声信号とに基づいて、反響音をキャンセルする分離行列である反響音キャンセル分離行列を算出し、当該音声取得部によって取得された当該残響音を含む既知の音声信号と当該残響音を含まない既知の音声信号と算出した前記反響音キャンセル分離行列とに基づいてブラインド音源分離における分離行列であるブラインド分離行列とブラインド残響除去における除去行列であるブラインド残響除去行列とを算出する残響データ演算部と、

前記残響データ演算部によって算出された前記反響音キャンセル分離行列に基づき残響抑圧を行うフィルタのフィルタ長を推定するフィルタ長推定部と、

前記フィルタ長推定部によって推定されたフィルタ長と前記残響データ演算部によって算出された反響音キャンセル分離行列と前記ブラインド分離行列と前記ブラインド残響除去行列とに基づき残響抑圧を行う残響抑圧部と、

を備えることを特徴とする残響抑圧装置。

【請求項2】

前記音声取得部によって取得された残響音を含む音声信号を周波数領域の第1信号に変換し、前記残響音を含まない既知の音声信号を周波数領域の第2信号に変換する短時間フ

ーリエ解析部

を備え、

前記残響データ演算部は、

前記短時間フーリエ解析部から入力された前記第 1 信号を周波数毎に空間球面化を行って第 3 信号を算出し、前記短時間フーリエ解析部から入力された前記第 2 信号を周波数毎にスケールの正規化を行うことで第 4 信号を算出し、算出された前記第 3 信号と前記第 4 信号と前記反響音キャンセル分離行列とを用いて前記ブラインド分離行列と前記ブラインド残響除去行列とを算出し、

前記残響抑圧部は、

前記フィルタ長推定部によって推定された前記フィルタ長に基づき、前記第 3 信号と前記第 4 信号と前記ブラインド分離行列と前記ブラインド残響除去行列と前記反響音キャンセル分離行列とを用いて、繰り返し独立成分分析処理を行い、独立成分分析処理が行われた演算結果に対してスケールリングを行い、スケールリングを行った信号からパワーが最大のものを選択することで前記音声取得部が取得した前記残響音を含む音声信号から前記残響音を抑圧する残響抑圧を行う

ことを特徴とする請求項 1 に記載の残響抑圧装置。

【請求項 3】

前記算出された前記反響音キャンセル分離行列、前記ブラインド分離行列、及び前記ブラインド残響除去行列に基づき残響時間を推定する残響特性推定部、を備え、

前記フィルタ長推定部は、

前記推定された残響時間に基づき前記フィルタ長を推定する

ことを特徴とする請求項 1 または請求項 2 に記載の残響抑圧装置。

【請求項 4】

前記フィルタ長推定部は、

直接音と間接音との比率に基づき前記フィルタ長を推定する

ことを特徴とする請求項 1 または請求項 2 に記載の残響抑圧装置。

【請求項 5】

当該残響抑圧装置の位置が変化したことを検出する検出部、

を更に備え、

残響データ演算部は、

前記残響抑圧装置の位置が変化したことを検出した場合に前記反響音キャンセル分離行列、前記ブラインド分離行列、及び前記ブラインド残響除去行列を演算する

ことを特徴とする請求項 1 から請求項 4 のいずれか 1 項に記載の残響抑圧装置。

【請求項 6】

前記検出部は、

前記残響抑圧装置の位置が変化したことを検出した場合に、前記残響抑圧部が残響抑圧に用いるパラメータ、あるいは、前記フィルタ長推定部がフィルタ長推定に用いるパラメータの少なくとも一方のパラメータを前記残響抑圧装置の位置が変化したことに基づき切り替える

ことを特徴とする請求項 5 に記載の残響抑圧装置。

【請求項 7】

テスト音声信号を出力する音声出力部、

を更に備え、

前記音声取得部は、前記出力されたテスト音声信号を取得し、残響データ演算部は、前記取得されたテスト音声信号から前記反響音キャンセル分離行列、前記ブラインド分離行列、及び前記ブラインド残響除去行列を演算する

ことを特徴とする請求項 1 から請求項 6 のいずれか 1 項に記載の残響抑圧装置。

【請求項 8】

残響抑圧装置の残響抑圧方法において、

音声取得部が、音声信号を取得する音声取得工程と、

10

20

30

40

50

残響データ演算部が、テスト音声信号を出力する音声出力部から既知の音声信号を出力させたときに、前記音声取得工程によって取得された残響音を含む前記既知の音声信号と残響音を含まない当該既知の音声信号とに基づいて、反響音をキャンセルする分離行列である反響音キャンセル分離行列を算出し、当該音声取得部によって取得された当該残響音を含む既知の音声信号と当該残響音を含まない既知の音声信号と算出した前記反響音キャンセル分離行列とに基づいてブライント音源分離における分離行列であるブライント分離行列とブライント残響除去における除去行列であるブライント残響除去行列とを算出する残響データ演算工程と、

フィルタ長推定部が、前記残響データ演算工程によって算出された前記反響音キャンセル分離行列に基づき残響抑圧を行うフィルタのフィルタ長を推定するフィルタ長推定工程と、

残響抑圧部が、前記フィルタ長推定工程によって推定されたフィルタ長と前記残響データ演算工程によって算出された反響音キャンセル分離行列と前記ブライント分離行列と前記ブライント残響除去行列とに基づき残響抑圧を行う残響抑圧工程と、

を備えることを特徴とする残響抑圧方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、残響抑圧装置、及び残響抑圧方法に関する。

【背景技術】

【0002】

残響抑圧処理は、遠隔会議通話または補聴器における明瞭度の向上およびロボットの音声認識（ロボット聴覚）に用いられる自動音声認識の認識率の向上を目的として、自動音声認識の前処理として利用されている重要な技術である。残響抑圧処理において、所定のフレーム毎に、取得した音声信号から残響成分を算出し、取得した音声信号から算出した残響成分を除去することで残響を抑圧していた（例えば、特許文献1参照）。

【先行技術文献】

【特許文献】

【0003】

【特許文献1】特開平9 261133号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、特許文献1の従来技術では、所定のフレーム長さにおいて残響抑圧を行っていたため、フレーム長が長い場合は処理に時間がかかりすぎるという問題点があり、フレーム長が短すぎると十分な残響抑圧の効果が得られにくいという問題点があった。

【0005】

本発明は、上記の問題点に鑑みてなされたものであって、精度良く残響抑圧を行える残響抑圧装置及び残響抑圧方法を提供することを課題としている。

【課題を解決するための手段】

【0006】

上記目的を達成するため、本発明に係る残響抑圧装置は、音声信号を取得する音声取得部と、音声信号を取得する音声取得部と、音声信号を出力する音声出力部と、前記音声出力部から既知の音声信号を出力させたときに、前記音声取得部によって取得された残響音を含む前記既知の音声信号と残響音を含まない当該既知の音声信号とに基づいて、反響音をキャンセルする分離行列である反響音キャンセル分離行列を算出し、当該音声取得部によって取得された当該残響音を含む既知の音声信号と当該残響音を含まない既知の音声信号と算出した前記反響音キャンセル分離行列とに基づいてブライント音源分離における分離行列であるブライント分離行列とブライント残響除去における除去行列であるブライント残響除去行列とを算出する残響データ演算部と、前記残響データ演算部によって算出さ

10

20

30

40

50

れた前記反響音キャンセル分離行列に基づき残響抑圧を行うフィルタのフィルタ長を推定するフィルタ長推定部と、前記フィルタ長推定部によって推定されたフィルタ長と前記残響データ演算部によって算出された反響音キャンセル分離行列と前記ブラインド分離行列と前記ブラインド残響除去行列とに基づき残響抑圧を行う残響抑圧部と、を備えることを特徴としている。

また、本発明に係る残響抑圧装置は、前記音声取得部によって取得された残響音を含む音声信号を周波数領域の第1信号に変換し、前記残響音を含まない既知の音声信号を周波数領域の第2信号に変換する短時間フーリエ解析部を備え、前記残響データ演算部は、前記短時間フーリエ解析部から入力された前記第1信号を周波数毎に空間球面化を行って第3信号を算出し、前記短時間フーリエ解析部から入力された前記第2信号を周波数毎にスケールの正規化を行うことで第4信号を算出し、算出された前記第3信号と前記第4信号と前記反響音キャンセル分離行列とを用いて前記ブラインド分離行列と前記ブラインド残響除去行列とを算出し、前記残響抑圧部は、前記フィルタ長推定部によって推定された前記フィルタ長に基づき、前記第3信号と前記第4信号と前記ブラインド分離行列と前記ブラインド残響除去行列と前記反響音キャンセル分離行列とを用いて、繰り返し独立成分分析処理を行い、独立成分分析処理が行われた演算結果に対してスケーリングを行い、スケーリングを行った信号からパワーが最大のものを選択することで前記音声取得部が取得した前記残響音を含む音声信号から前記残響音を抑圧する残響抑圧を行うようにしてもよい。

【0007】

また、本発明に係る残響抑圧装置において、前記算出された前記反響音キャンセル分離行列、前記ブラインド分離行列、及び前記ブラインド残響除去行列に基づき残響時間を推定する残響特性推定部、を備え、前記フィルタ長推定部は、前記推定された残響時間に基づき前記フィルタ長を推定するようにしてもよい。

【0008】

また、本発明に係る残響抑圧装置において、前記フィルタ長推定部は、直接音と間接音との比率に基づき前記フィルタ長を推定するようにしてもよい。

【0009】

また、本発明に係る残響抑圧装置において、当該残響抑圧装置の位置が変化したことを検出する検出部、を更に備え、残響データ演算部は、前記残響抑圧装置の位置が変化したことを検出した場合に前記反響音キャンセル分離行列、前記ブラインド分離行列、及び前記ブラインド残響除去行列を演算するようにしてもよい。

【0010】

また、本発明に係る残響抑圧装置において、前記検出部は、前記残響抑圧装置の位置が変化したことを検出した場合に、前記残響抑圧部が残響抑圧に用いるパラメータ、あるいは、前記フィルタ長推定部がフィルタ長推定に用いるパラメータの少なくとも一方のパラメータを前記残響抑圧装置の位置が変化したことに基づき切り替えるようにしてもよい。

【0011】

また、本発明に係る残響抑圧装置において、テスト音声信号を出力する音声出力部、を更に備え、前記音声取得部は、前記出力されたテスト音声信号を取得し、残響データ演算部は、前記取得されたテスト音声信号から前記反響音キャンセル分離行列、前記ブラインド分離行列、及び前記ブラインド残響除去行列を演算するようにしてもよい。

【0012】

上記目的を達成するため、本発明に係る残響抑圧装置における残響抑圧方法は残響抑圧装置の残響抑圧方法において、音声取得部が、音声信号を取得する音声取得工程と、残響データ演算部が、テスト音声信号を出力する音声出力部から既知の音声信号を出力させたときに、前記音声取得工程によって取得された残響音を含む前記既知の音声信号と残響音を含まない当該既知の音声信号とに基づいて、反響音をキャンセルする分離行列である反響音キャンセル分離行列を算出し、当該音声取得部によって取得された当該残響音を含む既知の音声信号と当該残響音を含まない既知の音声信号と算出した前記反響音キャンセル

分離行列とに基づいてブライント音源分離における分離行列であるブライント分離行列とブライント残響除去における除去行列であるブライント残響除去行列とを算出する残響データ演算工程と、フィルタ長推定部が、前記残響データ演算工程によって算出された前記反響音キャンセル分離行列に基づき残響抑圧を行うフィルタのフィルタ長を推定するフィルタ長推定工程と、残響抑圧部が、前記フィルタ長推定工程によって推定されたフィルタ長と前記残響データ演算工程によって算出された反響音キャンセル分離行列と前記ブライント分離行列と前記ブライント残響除去行列とに基づき残響抑圧を行う残響抑圧工程と、を備えることを特徴としている。

【発明の効果】

【0013】

本発明によれば、取得された音声信号から残響データを演算して、演算された残響データに基づいて残響特性を推定して、推定された残響特性に基づいて残響抑圧を行うフィルタのフィルタ長を推定するため、残響特性に応じた残響抑圧を精度良く効率的に行うことが可能になる。

【0014】

本発明によれば、推定された残響特性の残響時間に基づいてフィルタ長を推定するようにしたので、さらに精度が良く効率の良い残響抑圧を行うことが可能になる。

【0015】

本発明によれば、直接音と反射音との比率に基づいてフィルタ長を推定するようにしたので、さらに精度が良く効率の良い残響抑圧を行うことが可能になる。

【0016】

本発明によれば、当該残響抑圧装置が設置されている位置が変化したか否かを検出し、設置位置が変化して設置されている環境が変化した場合、残響データの演算と残響特性の推定を行い、推定された残響特性に基づいて残響抑圧を行うフィルタのフィルタ長を推定するため、さらに精度が良く効率の良い残響抑圧を行うことが可能になる。

【0017】

本発明によれば、残響抑圧部が残響抑圧に用いるパラメータ、あるいは、フィルタ長を推定するためのパラメータの少なくともどちらか一方のパラメータを予め設定されている位置に関する情報に基づいて切り替えるため、さらに精度が良く効率の良い残響抑圧を行うことが可能になる。

【0018】

本発明によれば、音声出力部が残響データを演算するためのテスト音声信号を出力して、音声取得部が、出力されたテスト音声信号を取得して、取得された音声信号から残響データを演算して、演算された残響データに基づいて残響特性を推定して、推定された残響特性に基づいて残響抑圧を行うフィルタのフィルタ長を推定するため、さらに精度が良く効率の良い残響抑圧を行うことが可能になる。

【図面の簡単な説明】

【0019】

【図1】本実施形態に係る残響抑圧装置を組み込んだロボットが取得する音声信号の一例を説明する図である。

【図2】同実施形態に係る残響抑圧装置100のブロック図の一例を示す図である。

【図3】同実施形態に係るSTFT処理を説明する図である。

【図4】同実施形態に係るMCSE-ICA部114の内部構成を説明する図である。

【図5】同実施形態に係る残響強度を検出する処理手順を説明する図である。

【図6】同実施形態に係るロボットのみが発話してマイクから音声信号を取得している状態を説明する図である。

【図7】同実施形態に係る残響強度の一例を示す図である。

【図8】同実施形態に係るMCSE-IC処理の変化の一例を示す図である。

【図9】同実施形態に係る実験に用いたデータ及び残響抑圧装置の設定条件である。

【図10】同実施形態に係る音声認識の設定を説明する図である。

10

20

30

40

50

【図 1 1】同実施形態に係る音声認識の設定を説明する図である。

【図 1 2】同実施形態に係る推定されたフィルタ長を用いた音声認識率の一例を示す図である。

【図 1 3】同実施形態に係るケース B（バージ・インの発生なし）且つ場所 1 の場合の音声認識率を示すグラフである。

【図 1 4】同実施形態に係るケース B（バージ・インの発生なし）且つ場所 2 の場合の音声認識率を示すグラフである。

【図 1 5】同実施形態に係るケース C（バージ・インの発生あり）且つ場所 1 の場合の音声認識率を示すグラフである。

【図 1 6】同実施形態に係るケース C（バージ・インの発生あり）且つ場所 2 の場合の音声認識率を示すグラフである。

10

【図 1 7】第 2 実施形態に係る残響抑圧装置 1 0 0 a のブロック図の一例を示す図である。

【発明を実施するための形態】

【0 0 2 0】

以下、図 1 ~ 図 1 7 を用いて本発明の実施形態について詳細に説明する。なお、本発明は斯かる実施形態に限定されず、その技術思想の範囲内で種々の変更が可能である。

【0 0 2 1】

[ 第 1 実施形態 ]

図 1 は、本実施形態における残響抑圧装置を組み込んだロボットが取得する音声信号の一例を説明する図である。ロボット 1 は、図 1 に示すように、基体部 1 1 と、基体部 1 1 にそれぞれ可動連結される頭部 1 2（可動部）と、脚部 1 3（可動部）と、腕部 1 4（可動部）とを備えている。また、ロボット 1 は、背負う格好で基体部 1 1 に収納部 1 5 を装着している。なお、基体部 1 1 には、スピーカ 2 0（音声出力部 1 4 0）が収納され、頭部 1 2 にはマイク 3 0 が収納されている。なお、図 1 は、ロボット 1 を側面から見た図であり、マイク 3 0 およびスピーカ 2 0 はそれぞれ複数収納されている。

20

【0 0 2 2】

まず、本実施形態の概略を説明する。

図 1 のように、ロボット 1 のスピーカ 2 0 から出力される音声信号を、ロボット 1 の発話  $S_r$  として説明する。

30

ロボット 1 が発話している時に、ヒト 2 が割り込んで発話することをバージ・イン（B a r g e - i n）と呼ぶ。バージ・インが発生しているとき、ロボット 1 には、当該ロボット 1 の発話のために、割り込んできたヒト 2 の発話を聞き分けることが困難である。

そして、ヒト 2 およびロボット 1 が発話している場合、ロボット 1 のマイク 3 0 には、ヒト 2 の発話  $S_u$  が空間を經由して伝達する残響音を含むヒト 2 の音声信号  $h_u$  と、ロボット 1 の発話  $S_r$  が空間を經由して伝達する残響音を含むロボット 1 の音声信号  $h_r$  とが入力される。

【0 0 2 3】

図 1 において、ロボット 1 のマイク 3 0 が集音する音声信号をモデル化すると、 $h_u + h_r = H_u \cdot S_u + H \cdot S_r$  のように表せる。 $H_u$  と  $H$  は周波数領域の関数である。 $H_u \cdot S_u + H \cdot S_r$  において、 $S_r$  はロボット 1 の発話のため、当該ロボット 1 にとって既知である。マイク 3 0 が集音した音声信号の中で  $H_u \cdot S_u$  には、ヒト 2 が発話してからロボット 1 に伝播する間に残響音（エコー）が付加されてしまっているため、 $H_u \cdot S_u$  を用いて音声認識するより、 $S_u$  を用いて音声認識を行えば認識率が高いことが予測される。また、 $H$  は、ロボット 1 が単独でスピーカ 2 0 を介して発話し、発話した音声データを、マイク 3 0 を介して取得し、当該ロボット 1 がいる環境の残響特性を解析することで算出する。さらに、本実施形態では、ICA（i n d e p e n d e n t c o m p o n e n t a n a l y s i s；独立成分分析）をベースにした MCSB-ICA（m u l t i - c h a n n e l s e m i - b l i n d I C A）を用いて残響音をキャンセル、すなわち抑圧する。さらに、MCSB-ICA の分離フィルタのフレーム数を、算出した残

40

50

響特性に基づいて推定することで、ロボット1がいる環境に合わせたフレーム数を算出する。そして、最終的には、算出されたフレーム数を用いて残響成分を抑圧することでヒト2の発話の音声信号 $S_{\text{ヒト2}}$ を算出する。

【0024】

図2は、本実施形態における残響抑圧装置100のブロック図の一例を示す図である。図2のように、残響抑圧装置100にはマイク30、スピーカ20が接続され、マイク30は複数のマイク31、32・・・を備えている。また、残響抑圧装置100は、制御部101と、音声生成部102と、音声出力部103と、音声取得部111と、残響データ算出部112と、STFT部113と、MCSCB-ICA部114と、記憶部115と、フィルタ長推定部116と、分離データ出力部117とを備えている。

10

【0025】

制御部101は、残響特性を測定するための音声を生成して出力する指示を音声生成部102に出力し、ロボット1が残響特性を測定するための発話中を示す信号を音声取得部111とMCSCB-ICA部114に出力する。

【0026】

音声生成部102は、制御部101からの指示に基づき、残響特性測定用の音声信号(テスト信号)を生成し、生成した音声信号を音声出力部103に出力する。

【0027】

音声出力部103には、生成された音声信号が入力され、入力された音声信号を所定のレベルに増幅してスピーカ20に出力する。

20

【0028】

音声取得部111は、マイク30が集音した音声信号を取得し、取得した音声信号をSTFT部113に出力する。また、音声取得部111は、制御部101から残響特性を測定するための音声を生成して出力する指示が入力された時、残響特性を測定するための音声信号を取得し、取得した音声信号を残響データ算出部112に出力する。

【0029】

残響データ算出部(残響データ演算部)112には、取得された音声信号と生成された音声信号が入力され、取得された音声信号と生成された音声信号、および記憶部115に記憶されている演算式を用いて反響音キャンセル分離行列 $W_r$ を算出する。また、残響データ算出部112には、算出した反響音キャンセル分離行列 $W_r$ を記憶部115に書き込んで記憶させる。

30

【0030】

STFT(short-time Fourier transformation; 短時間フーリエ解析)部113には、取得された音声信号と生成された音声信号が入力され、入力された各音声信号にハニング等の窓関数を音声信号に乗じて有限期間内で、解析位置をシフトしながら解析を行う。そして、STFT部113は、取得された音声信号を、フレーム $t$ 毎にSTFT処理して時間-周波数領域の信号 $x(\quad, t)$ に変換し、また、生成された音声信号を、フレーム $t$ 毎にSTFT処理して時間-周波数領域の信号 $s_r(\quad, t)$ に変換し、変換した信号 $x(\quad, t)$ と信号 $s_r(\quad, t)$ を周波数ごとにMCSCB-ICA部114に出力する。図3(a)と図3(b)は、STFT処理を説明する図である。図3(a)は、取得された音声信号の波形であり、図3(b)は、この取得された音声信号に乗じられる窓関数である。図3(b)において、記号 $U$ はシフト長であり、記号 $T$ は解析を行う期間(窓長)を示している。

40

【0031】

MCSCB-ICA部(残響抑圧部)114には、STFT部113から変換された信号 $x(\quad, t)$ と信号 $s_r(\quad, t)$ が周波数ごとに入力され、制御部101からロボット1が残響特性を測定するための発話中を示す信号が入力され、フィルタ長推定部116から推定されたフィルタ長データが入力される。また、MCSCB-ICA部114は、残響特性を測定するための発話中を示す信号が入力されていない場合、入力された信号 $x(\quad, t)$ と信号 $s_r(\quad, t)$ と記憶部114に記憶されている反響音キャンセル分離行

50

列  $W_r$ 、各モデル及び各係数を用いて、分離フィルタ  $W_{1u}$  と  $W_{2u}$  を算出する。分離フィルタ  $W_{1u}$  と  $W_{2u}$  算出後、マイク 30 が取得した音声信号からヒト 2 の直接発話信号を分離し、分離した直接発話信号を分離データ出力部 117 に出力する。

【0032】

図 4 は、MCSB-ICA 部 114 の内部構成を説明する図である。図 4 のように、STFT 部 113 から入力された信号  $x(\cdot, t)$  はバッファ 201 を介して強制空間球面化部 211 に入力され、STFT 部 113 から入力された信号  $s_r(\cdot, t)$  はバッファ 202 を介して分散正規化部 212 に入力される。そして、ICA 部 221 には、強制空間球面化部 211 から空間球面化された信号と、分散正規化部 212 から正規化された信号とが入力され、入力された信号を用いて繰り返し ICA 処理を行い、演算結果をスケールリング部 231 に出力し、スケールリングされた信号を直接発話分離部 241 に出力する。なお、スケールリング部 231 は、projection Back 処理を用いてスケールリングを行い、直接発話分離部 241 は、入力された信号からパワーが最大のものを選択して出力する。

10

【0033】

記憶部 115 には、ロボット 1 がマイク 30 を介して取得する音声信号のモデル、解析するための分離モデル、解析するために必要なパラメータ等が予め書き込まれて記憶され、さらに、算出された反響音キャンセル分離行列  $W_r$ 、分離フィルタ  $W_{1u}$  及び分離フィルタ  $W_{2u}$  が書き込まれて記憶されている。

【0034】

20

フィルタ長推定部（残響特性推定部、フィルタ長推定部）116 は、記憶部 115 に記憶されている反響音キャンセル分離行列  $W_r$  を読み出し、読み出した反響音キャンセル分離行列  $W_r$  から後述する方法でフィルタ長を推定し、推定したフィルタ長データを MCSB-ICA 部 114 に出力する。なお、フィルタ長とは、フレーム（窓）をサンプリングする数に関する値であり、フィルタ長が大きくなると時間方向に長い期間、サンプリングを行うことになる。

【0035】

分離データ出力部 117 には、MCSB-ICA 部 114 から分離された直接発話信号が入力され、入力された直接発話信号を、例えば非図示の音声認識部に出力する。

【0036】

30

次に、ロボット 1 が取得した音声から必要な音声信号を分離するための分離モデルについて説明する。記憶部 115 には、ロボット 1 がマイク 30 を介して取得する音声信号は、式 (1) の FIR (finite impulse response; 有限インパルス応答) のモデルのように定義する。

【0037】

【数 1】

$$x(t) = \sum_{n=0}^N h_u(n) s_u(t-n) + \sum_{m=0}^M h_r(m) s_r(t-n) \quad \dots (1)$$

40

【0038】

式 (1) において、記号  $x_1(t) \dots x_L(t)$  は、各マイク 31 ~ 32 の各スペクトル ( $L$  はマイク番号)、 $x(t)$  はベクトルであり  $[x_1(t), x_2(t), \dots, x_L(t)]^T$ 、 $s_u(t)$  はヒト 2 の発話、 $s_r(t)$  は既知のロボット 1 のスペクトル、 $h_u(n)$  はヒト 2 の音声スペクトルの  $N$  次元の FIR 係数ベクトル、 $h_r(m)$  は既知のロボット 1 の  $M$  次元の FIR 係数ベクトルである。式 (1) は、ロボット 1 がマイク 30 を介して取得する時刻  $t$  におけるモデル化である。

【0039】

また、記憶部 115 には、ロボット 1 のマイク 30 が集音した音声信号について、式 (

50



2)のように残響成分を含んだベクトル $X(t)$ としてモデル化され予め記憶されている。さらに、記憶部115には、ロボット1の発話の音声信号について、式(3)のように残響成分を含んだベクトル $S_r(t)$ としてモデル化されて予め記憶されている。

【0040】

【数2】

$$X(t) = [x(t), x(t-1), \dots, x(t-N)]^T \quad \dots(2)$$

10

【0041】

【数3】

$$S_r(t) = [s_r(t), s_r(t-1), \dots, s_r(t-M)]^T \quad \dots(3)$$

【0042】

式(3)において、 $s_r(t)$ はロボット1が発話した音声信号であり、 $s_r(t-1)$ は空間を伝達されて「1」遅延して音声信号が届くことを表し、 $s_r(t-M)$ は「M」遅延して届く音声信号が届くことを表している。すなわち、ロボット1から離れている距離が大きく、遅延量が大きいくほど残響成分が大きくなることを表している。

20

【0043】

次に、ICAを用いて既知の直接音 $S_r(t)$ と $X(t-d)$ と、ヒト2の直接発話信号 $s_u$ とを独立となるように分離するため、記憶部115には、MCSB-ICAの分離モデルが次式(4)のように定義し、記憶部115に記憶されている。

【0044】

【数4】

$$\begin{pmatrix} \hat{s}(t) \\ X(t-d) \\ S_r(t) \end{pmatrix} = \begin{pmatrix} W_{1u} & W_{2u} & W_r \\ 0 & I_2 & 0 \\ 0 & 0 & I_r \end{pmatrix} \begin{pmatrix} x(t) \\ X(t-d) \\ S_r(t) \end{pmatrix} \quad \dots(4)$$

30

【0045】

式(4)において、 $d$ (0より大きい)は、初期反射間隔であり、 $X(t-d)$ は、 $X(t)$ を $d$ 遅延させたベクトルであり、式(5)は、 $L$ 次元の推定された信号ベクトルである。

【0046】

【数5】

$$\hat{s}(t) \quad \dots(5)$$

40

【0047】

また、 $W_{1u}$ は、 $L \times L$ のブラインド分離行列(分離フィルタ)、 $W_{2u}$ は、 $L \times L(N+1)$ のブラインド残響除去行列(分離フィルタ)、 $W_r$ は、 $L \times (M+1)$ の残響音キャンセルの分離行列(取得した残響特性に基づく残響要素)である。

また、 $I_2$ と $I_r$ は、それぞれに対応した大きさの単位行列である。そして、式(5)

50

には、ヒト2の発話の直接発話信号といくつかの反射音信号とを含んでいる。

【0048】

次に、式(4)を解くためのパラメータについて説明する。式(4)において、分離パラメータのセット  $W = \{W_{1u}, W_{2u}, W_r\}$  を、結合確率密度関数 (probability density function) と  $s(t)$ 、 $X(t-d)$  および  $S_r(t)$  の周辺確率密度関数 (個々のパラメータの独立な確率分布を表わす周辺確率密度関数) の積との間の差の尺度として KL (kullback-Leibler; カルバック・ライブラー) 情報量を最小化するように推定する。また、周波数における分離行列の初期値  $W_{1u}(\cdot)$  は、周波数  $\cdot + 1$  において推定行列  $W_{1u}(\cdot + 1)$  にセットされている。

10

【0049】

MCSB-ICA部114は、分離パラメータのセット  $W$  を、KL情報量を自然勾配法により最小にするように各分離フィルタ次式(6)~式(9)のルールに従い繰り返し更新することで推定を行う。また、式(6)~式(9)は、記憶部115に予め書き込まれて記憶されている。

【0050】

【数6】

$$D = \Lambda - E[\phi(\hat{s}(t))\hat{s}^H(t)] \quad \dots(6)$$

20

【0051】

【数7】

$$W_{1u}^{[j+1]} = W_{1u}^{[j]} + \mu DW_{1u}^{[j]} \quad \dots(7)$$

【0052】

【数8】

$$W_{2u}^{[j+1]} = W_{2u}^{[j]} + \mu(DW_{2u}^{[j]} - E[\phi(\hat{s}(t))X^H(t-d)]) \quad \dots(8)$$

30

【0053】

【数9】

$$W_r^{[j+1]} = W_r^{[j]} + \mu(DW_r^{[j]} - E[\phi(\hat{s}(t))S_r^H(t)]) \quad \dots(9)$$

【0054】

なお、式(6)、式(8)~式(9)において、上付きHは共役転置演算(エルミート転置)を表す。また、式(5)において、 $\Lambda$  は非ホロノミック拘束行列、すなわち、式(10)の対角行列である。

40

【0055】

【数10】

$$E[\phi(\hat{s}(t))\hat{s}^H(t)] \quad \dots(10)$$

50

## 【 0 0 5 6 】

また、式(7)～式(9)において、 $u$ は、ステップ・サイズのパラメータであり、 $(x)$ は、非線形関数ベクトル $[(x_1), \dots, (x_L)]^H$ であり、次式(11)のように表され、記憶部115に書き込まれて記憶されている。

## 【 0 0 5 7 】

## 【 数 1 1 】

$$\phi(x) = -\frac{d}{dx} \log p(x) \quad \dots(11)$$

10

## 【 0 0 5 8 】

さらに、音源のPDFは、分散量 $\sigma^2$ であるとした場合、雑音に強いPDFである $p(x) = \exp(-|x|^2 / \sigma^2) / (2\sigma^2)$ であり、 $(x) = x^* / (2\sigma^2 |x|)$ であり、 $x^*$ は $x$ の共役であると仮定する。この2つの関数は、連続領域である $|x| > 0$ において定義される。

## 【 0 0 5 9 】

次に、音声を分離する処理手順を、図5～図8を用いて説明する。図5は、本実施形態における残響強度を検出する処理手順を説明する図である。なお、残響強度の検出は、ロボット1がいる環境が変わった場合、例えば、別の部屋に移動した後、室外に出た後毎に行う。また、ロボット1は、例えば、当該ロボット1に組み込まれている非図示のカメラで撮像された画像データを用いて、環境が変化したか否かを判定する。あるいは、ロボット1が水平方向または垂直方向に移動し、当該ロボット1がいた位置が変化した場合にも残響強度を検出する処理を行うようにしてもよい。

20

## 【 0 0 6 0 】

## [ステップS1; Emission of self speech]

まず、図6のように、ロボット1は、当該ロボット1が現在いる環境で、制御部101は、残響強度を測定するための所定の音声信号を生成する指示を音声生成部102に出力する。音声生成部102には、所定の音声信号を生成する指示が入力され、入力された生成指示に基づき所定の音声信号を生成し、生成した所定の音声信号を音声出力部103に出力する。音声出力部103には、生成された所定の音声信号が入力され、入力された所定の音声信号を所定のレベルに増幅してスピーカ20に出力する。なお、残響強度を測定するための所定の音声信号は、例えば、1つの母音または1つの子音であってもよい。図6は、ロボットのみが発話してマイクから音声信号を取得している状態を説明する図である。

30

## 【 0 0 6 1 】

次に、音声取得部111には、マイク30が集音した音声信号が入力され、入力された音声信号を残響データ算出部112に出力する。マイク30が集音する音声信号は、音声生成部102が生成した音声信号 $S_r$ に、スピーカ20から発せられた音声 $h_r$ が壁、天井、床などで反響した残響成分を含む音声信号 $h_r$ である。

## 【 0 0 6 2 】

次に、残響データ算出部112には、取得された音声信号が入力され、入力された音声信号を記憶部115に記憶されている式(9)を用いて反響音キャンセル分離行列 $W_r$ を算出する。また、残響データ算出部112は、演算した残響特性データを記憶部115に書き込んで記憶させる。なお、式(9)を演算するとき、入力値は $W_r$ のみなのでフィルタ長を1に設定する。

40

## 【 0 0 6 3 】

## [ステップS2; Calculation of echo intensities]

ステップS2では、ステップS1で算出された $W_r$ を使って、フィルタ長を推定するための残響強度のグラフを生成する。

まず、フィルタ長推定部116は、記憶部115に記憶されている反響音キャンセル分

50

離行列  $W_r$  を読み出す。フィルタ長推定部 116 は、読み出した反響音キャンセル分離行列  $W_r$  を、パラメータ  $W_r$  を式 (12) のような行列に置き直す。

【0064】

$W_r = [w_r(0) w_r(1) \cdots w_r(M)] \cdots (12)$

【0065】

なお、式 (12) の  $W_r$  において、 $w_r(m)$  は、 $L \times 1$  ベクトルであり式 (13) のように表される。

【0066】

【数12】

$$W_r(m) = [w_r^1(m) w_r^2(m) \cdots w_r^L(M)]^T \cdots (13)$$

10

【0067】

そして、周波数  $\omega$  におけるこのフィルタの正規化されたパワー関数は、次式 (14) のように定義する。

【0068】

【数13】

$$p_r^i(\omega, m) = \frac{|\omega_r^i(\omega, m)|^2}{\max_m |\omega_r^i(\omega, m)|^2} \cdots (14)$$

20

【0069】

式 (14) において、 $i$  はマイク 30 の番号 (マイク 31、32、 $\cdots$ ) であり、 $m$  はフィルタのインデックスである。式 (14) のパワー関数は、残響強度を反映し、また、環境の残響時間に関係しているため、このパワー関数に基づいて残響時間を推定する。

30

次に、平均化された周波数のパワー関数と平均化されたマイクのパワー関数  $P$  と、関数  $P$  の対数値  $L$  は、次式 (15) と式 (16) のように残響時間のための基準として定義する。

【0070】

【数14】

$$p(m) = \frac{\sum_i \sum_{\omega \in \Omega} p_r^i(\omega, m)}{\max_m \sum_i \sum_{\omega \in \Omega} p_r^i(\omega, m)} \cdots (15)$$

40

【0071】

【数15】

$$L(m) = 20 \log_{10} P(m) \cdots (16)$$

【0072】

50

式(15)において、 $L_d$ は周波数バンド・セットに基づく値である。フィルタ長推定部116は、この式(15)と式(16)を用いて、図7のように残響強度を仮想的にプロットする。図7において、縦軸は音声レベルであり、横軸は時間軸を表している。図7のように、生成された音声信号をスピーカ30から発した時(時刻0)の音声レベルが一番高くロボット1がいる環境の残響特性に応じて、音声レベルは下がっていく。

【0073】

[ステップS3; Estimation of dereverberation filter length]

ステップS3では、図7のプロットされた残響強度のグラフを用いて、フィルタ長Mを検定する。

10

まず、図7のように、フィルタ長推定部116は、フィルタ長の推定のため式(17)を用いて線形回帰解析を行う。

【0074】

$$y = a \times m + b \quad \dots (17)$$

【0075】

式(17)において、 $a$ と $b$ は係数であり、 $m$ はフィルタ長のインデックス、そして $y$ は $L(m)$ と等価である。次に、図7のように、フィルタ長推定部116は、 $P(m)$ のピーク値からいくつかのサンプルを抽出し、最小二乗平均(LMS; least mean square)法を用いて $a$ と $b$ を推定する。

次に、フィルタ長推定部116は、残響除去のフィルタ長を、次式(18)において、 $m$ が $L(m) = L_d$ の値を満足するように算出し、算出した残響除去のフィルタ長をICA部221に出力する。

20

【0076】

【数16】

$$\hat{N} = \frac{L_d - b}{a} \quad \dots (18)$$

30

【0077】

一例として、図7において、 $RT_{20} = 240 \text{ msec}$  ( $RT_{20}$ は残響時間)、そして線形回帰線251を式(17)により推定する。そして、推定されたフィルタ長は、式(18)において $L_d = -60$  (ライン252)との交点253の値、 $M = \text{約} 13$ である。

【0078】

[ステップS4; Incremental separation poling notification]

ヒト2の発話が発声した場合、このステップS4を行い、式(4)を用いて式(5)を求め、マイク30から取得された音声信号からヒト2の残響成分除去した音声信号を算出する。

40

【0079】

音声取得部111には、マイク30が集音した音声信号が入力され、入力された音声信号をSTFT部113に出力する。また、音声生成部102は、音声を生成している場合、生成した音声信号をSTFT部113に出力する。

【0080】

次に、STFT部113には、マイク30が取得した音声信号と、音声生成部102が生成した音声信号とが入力され、取得された音声信号をフレーム $t$ 毎にSTFT処理して時間-周波数領域の信号 $x(\cdot, t)$ に変換し、変換した信号 $x(\cdot, t)$ を周波数ごとにMCMB-ICA部114に出力する。また、STFT部113は、生成された音声

50

信号を、フレーム  $t$  毎に STFT 処理して時間 - 周波数領域の信号  $s_r(\cdot, t)$  に変換し、変換した信号  $s_r(\cdot, t)$  を周波数ごとに MCSB - ICA 部 114 に出力する。

【0081】

MCSB - ICA 部 114 の強制空間球面化部 211 には、変換された信号  $x(\cdot, t)$  が周波数ごとに入力され、周波数をインデックスとして順次、次式(19)を用いて空間球面化を行い、 $z(t)$  を算出する。また、式(19)と式(21)は、式(5)を解く上で高速化を行うために用いている。

【0082】

【数17】

$$z(t) = V_u x(t) \quad \dots(19)$$

【0083】

ただし、 $V_u$  は式(20)である。

【0084】

【数18】

$$V_u = E_u \Lambda^{-\frac{1}{2}} E_u^H \quad \dots(20)$$

【0085】

さらに、式(20)において、 $E_u$  と  $\Lambda$  は、固有ベクトル行列であり、固有対角行列  $R_u = E | x(t) x^H(t) |$  である。

さらに、MCSB - ICA 部 114 の分散正規化部 212 には、変換された信号  $s_r(\cdot, t)$  が周波数ごとに入力され、周波数をインデックスとして順次、次式(21)を用いてスケールの正規化を行う。

【0086】

【数19】

$$\tilde{s}_r(t) = \lambda_r^{-\frac{1}{2}} s_r(t) \quad \dots(21)$$

【0087】

なお、スケージングの正規化において、逆変換法 (projection back method) を用い、逆分離行列の要素は、分離信号に従って乗算される。そして、式(22)の  $i$  番目の列、 $j$  番目の行の要素  $c_j$  は、式(5)の  $j$  番目の要素のスケージングは、式(23) ~ 式(24)の式の関係に従って行う。

【0088】

【数20】

$$\hat{H}_u = (W_{1u} V_0)^{-1} \quad \dots(22)$$

【0089】

10

20

30

40

【数 2 1】

$$l_j = \arg \max_l |\hat{H}_u(l, j)| \quad \dots(23)$$

【0090】

【数 2 2】

$$c_j = \hat{H}_u(l_j, j) \quad \dots(24)$$

10

【0091】

強制空間球面化部 2 1 1 は、このように演算された  $z(\quad, t)$  を I C A 部 2 2 1 に出力し、分散正規化部 2 1 2 は、このように演算された式 ( 2 1 ) の値を I C A 部 2 2 1 に出力する。

【0092】

次に、I C A 部 2 2 1 には、演算された  $z(\quad, t)$  と式 ( 2 1 ) の値とが入力され、さらに、記憶部 1 1 5 に記憶されている分離モデル ( 分離フィルター ) を読み出す。次に、I C A 部 2 2 1 は、式 ( 4 )、式 ( 6 ) ~ 式 ( 9 ) の  $x$  に式 ( 1 9 ) を代入し、 $s_r$  に式 ( 2 1 ) を代入して、 $W_{1u}$  と  $W_{2u}$  を算出し、すでにステップ S 1 で算出された  $W_r$  を用いて、M C S B - I C A 部 1 1 4 が式 ( 5 ) のデータを算出する。

20

【0093】

図 8 は、M C S B - I C A 処理の変化の一例を示す図である。通常分離モードにおいて、M C S B - I C A のブロック幅増加分離を行う。I C A は、分離行列を安定して推測するために、所定の持続時間、データをバッファする。このようにバッファを使用するため、時間  $t$  の分離を行うため先行するブロックサイズ  $I_b$  を利用する。図 8 においては、シフト量  $I_s$  が増加する場合、遅れ時間も増加する。また、シフト量  $I_s$  が減少する場合、算出処理が増加する。このように、本実施形態では、オーバーラップ・パラメータ係数  $I_s$  を使用する。

30

【0094】

次に、本実施形態の残響抑圧装置を備えるロボット 1 で行った実験方法と実験結果の例を説明する。図 9 ~ 図 1 2 は、実験条件である。図 9 は、実験に用いたデータ及び残響抑圧装置の設定条件である。図 9 のように、インパルス応答は 1 6 K H z サンプル、残響時間は 2 4 0 m s と 6 7 0 m s、ロボット 1 とヒト 2 との距離は 1 . 5 m、ロボット 1 とヒト 2 の角度は 0 度、4 5 度、9 0 度、- 4 5 度、- 9 0 度、使用したマイク 3 0 の本数は 2 本 ( ロボット 1 の頭部に設置 )、S T F T 分析はハニング窓のサイズ 3 2 m s ( 5 1 2 ポイント ) かつシフト量 1 2 m s ( 1 9 2 ポイント )、入力信号データは [ - 1 . 0 1 . 0 ] に正規化されたものである。

40

【0095】

図 1 0 は、音声認識の設定を説明する図である。図 1 0 のように、テスト・セットは 2 0 0 の文章 ( 日本語 )、訓練セットは 2 0 0 人 ( それぞれ 1 5 0 の文章 )、音響モデルは P T M - t r i p h o n e、3 値の H M M ( 隠れマルコフモデル )、言語モデルは語彙サイズ 2 0 k、発話解析はハニング窓のサイズ 3 2 m s ( 5 1 2 ポイント )、シフト量 1 0 m s、特徴量は M F C C ( M e l - F r e q u e n c y C e p s t r m C o e f f i c i e n t ; スペクトル包絡 ) は 2 5 次 ( 1 2 次 + 1 2 次 + パワー ) である。また、他の S T F T 設定条件は、フレーム間隔係数  $d = 2$ 、反響キャンセルのフィルタ長  $N$  と通常分離モードの残響除去のフィルタ長  $M$  は同じ値、適応ステップ・サイズのための係数は予め設定され、推定されたフィルタ係数は、 $= \{ 5, 6, \dots, 200 \}$  かつ  $L$

50

$d = -60$ 、直線回帰のためのサンプル数は6に設定してある。また、音声認識エンジンは、公知のJulius (<http://julius.sourceforge.jp/>) を使用している。

【0096】

次に、実験結果を図11~図16に示す。図11は、推定されたフィルタ長の設定を示した図である。図11のように、ノイズあり且つ残響時間が240msの場合、ノイズあり且つ残響時間670msの場合、ノイズなし且つ残響時間が240msの場合、ノイズなし且つ残響時間670msの場合、各々について $M_{max}$ が20, 30, 50についての推定されたフィルタ長の平均値と偏差を示している。場所1 (Env. I) は、通常の部屋 (残響時間 $RT_{20} = 240\text{ms}$ )、場所2 (Env. II) は、ホールのような部屋 (残響時間 $RT_{20} = 670\text{ms}$ ) である。

10

【0097】

図12は、推定されたフィルタ長を用いた音声認識率の一例を示す図である。図12のように、ケースBは、バージ・インが発生していない場合、ケースCは、バージ・インが発生している場合、各々について音声分離無しでの認識率 (no proc)、ブロックサイズ $I_b$ が166 (2秒)、208 (2.5秒)、255 (3秒)、残響時間240msと670msの各音声認識率を示している。また、シフト量 $I_s$ は、ブロックサイズ $I_b$ の半分に設定されている。一例として、残響音がないクリーンな音声信号による認識率は、実験に用いた残響抑圧装置では約93%である。

【0098】

図12の結果をグラフにしたのが図13~図16である。図13は、ケースB (バージ・インの発生なし) 且つ場所1の場合の音声認識率を示すグラフであり、図14は、ケースB (バージ・インの発生なし) 且つ場所2の場合の音声認識率を示すグラフである。図15は、ケースC (バージ・インの発生あり) 且つ場所1の場合の音声認識率を示すグラフであり、図16は、ケースC (バージ・インの発生あり) 且つ場所2の場合の音声認識率を示すグラフである。各グラフの横軸はフィルタ長 (N) であり、縦軸は音声認識率 (%) である。

20

図13のように、残響時間が短い部屋 (場所1) 且つバージ・インが発生していない場合、推定されたフィルタ長 ( $N = 14$ ) 301より不適切なフィルタ長 ( $N = 35$ ) の方が認識率 (正答率) は低くかつブロックサイズ $I_b$ を変えると認識率の差が大きくなる。フィルタ長 ( $N = 35$ ) 302の場合はブロックサイズ $I_b$ により認識率に差が生じている。一方、残響時間が長い部屋 (場所2) 且つバージ・インが発生していない場合、推定されたフィルタ長 ( $N = 35$ ) で認識率は60%以上である。そして、図13と図14のように、残響時間が短い場合のフィルタ長は $N = 14$ で短く、残響時間が長い場合のフィルタ長は $N = 36$ で長い。このように、ロボット1が取得した環境の残響特性に基づき、適切なフィルタ長 (フレーム長) を推定することで、音声認識率を改善できる。

30

図15のように、残響時間が短い部屋 (場所1) 且つバージ・インが発生している場合、推定されたフィルタ長 ( $N = 14$ ) より不適切なフィルタ長 ( $N = 35$ ) の方が認識率 (正答率) は低くかつブロックサイズ $I_b$ を変えると認識率の差が大きくなる。一方、残響時間が長い部屋 (場所2) 且つバージ・インが発生している場合、推定されたフィルタ長 ( $N = 35$ ) で認識率は40%以上である。

40

【0099】

以上のように、残響特性に応じて分離フィルタ長であるフレーム長を設定するようになったので、音声認識率が向上し、さらに音声認識にかかる演算量も適切にすることが可能になる。

【0100】

また、本実施形態では、残響特性として残響時間を用いた例を説明したが、D値 (音声の明瞭さを表す値であり、直接音が到達してから0~50msまでのパワーと、0~音声が減衰するまでのパワーの比) を用いても良い。

【0101】

また、本実施形態では、残響特性の測定を制御部101から残響特性を測定するための

50



音声を生成して出力する指示が入力された時、残響特性を測定するための音声信号を取得して残響特性を測定する例を説明したが、音声取得部 1 1 1 は、音声生成部 1 0 2 が出力する生成された音声信号と比較しながら取得し、取得中にバージ・インが発生しているか否かを判別して、バージ・インが発生していないときに残響特性の測定用の音声信号を取得するようにしてもよい。

#### 【 0 1 0 2 】

##### [ 第 2 実施形態 ]

次に、第 2 実施形態について、図 1 7 を用いて説明する。図 1 7 は、本実施形態における残響抑圧装置 1 0 0 a のブロック図の一例を示す図である。第 1 実施形態では、ロボット 1 は、環境が変わった場合に、発話を行い、当該ロボット 1 がいる環境の残響特性を測定する例を説明した。残響特性の測定は、例えば、ロボット 1 が移動する部屋毎に例えばマークが設置され、設置されているマークをロボット 1 のカメラ 4 0 が撮像して公知の画像認識の手法を用いて、マークを検出して環境、例えば部屋を移動したことを検出した場合に行う。あるいは、ロボット 1 の記憶部 1 1 4 に予めマップを書き込んで記憶させておき、マップに基づき環境変化を検出した場合に行う。

#### 【 0 1 0 3 】

図 1 7 のように、本実施形態における残響抑圧装置 1 0 0 a は、画像取得部 3 0 1 と、環境変化検出部 3 0 2 とをさらに備えている。また、残響抑圧装置 1 0 0 a には、カメラ 4 0 が接続され、画像取得部 3 0 1 には、カメラにより撮像された画像信号が入力され、入力された画像信号を環境変化検出部 3 0 2 に出力する。環境変化検出部 3 0 2 は、入力された画像信号に基づき、残響抑圧装置 1 0 0 a が組み込まれているロボット 1 a がいる位置が変化したか否かを判定し、位置が変化したことを検出した場合、位置が変化したことを示す信号を制御部 1 0 1 a に出力する。制御部 1 0 1 a は、位置が変化したことを示す信号が入力された場合、音声生成部 1 0 2 に残響特性測定用の音声信号（テスト信号）を生成する指示を出力する。以下、第 1 実施形態と同様の処理を行う。

#### 【 0 1 0 4 】

また、各パラメータを環境毎に予め記憶部 1 1 5 a に書き込んで記憶させておき、マップ、マークとおのおの関連づけて記憶部 1 1 5 a に記憶させておく。

そして、ロボット 1 a が、環境が変わったことを検出した場合、制御部 1 0 1 a は、残響特性を測定するとともに、各パラメータのセットを記憶部 1 1 4 a から読み出して切り替えるようにしても良い。

#### 【 0 1 0 5 】

また、記憶部 1 1 5 a に残響データが記憶されていない環境で、残響測定を行い、測定された残響特性と関連付けて、その環境に基づくパラメータを算出して、算出したパラメータを関連づけて新たに記憶部 1 1 5 a に記憶させるようにしてもよい。

#### 【 0 1 0 6 】

また、例えば、各部屋にロボット 1 a へ位置に関する情報を送信する非図示の位置情報送信装置を設置し、ロボット 1 a はこの位置情報を受信した場合に環境が変化すると検出して、残響特性を測定するようにしてもよい。

#### 【 0 1 0 7 】

なお、第 1、第 2 実施形態では、残響抑圧装置 1 0 0 及び残響抑圧装置 1 0 0 a をロボット 1 ( 1 a ) に組み込んだ例を説明したが、残響抑圧装置 1 0 0 及び残響抑圧装置 1 0 0 a は、例えば音声認識装置、音声認識装置を有する装置などに組み込んで用いることも可能である。

#### 【 0 1 0 8 】

なお、実施形態の図 2 及び図 1 7 の各部の機能を実現するためのプログラムをコンピュータ読み取り可能な記録媒体に記録して、この記録媒体に記録されたプログラムをコンピュータシステムに読み込ませ、実行することにより各部の処理を行ってもよい。なお、ここでいう「コンピュータシステム」とは、OS や周辺機器等のハードウェアを含むものとする。

10

20

30

40

50

また、「コンピュータシステム」は、WWWシステムを利用している場合であれば、ホームページ提供環境（あるいは表示環境）も含むものとする。

また、「コンピュータ読み取り可能な記録媒体」とは、フレキシブルディスク、光磁気ディスク、ROM (Read Only Memory)、CD-ROM等の可搬媒体、USB (Universal Serial Bus) I/F (インタフェース) を介して接続されるUSBメモリー、コンピュータシステムに内蔵されるハードディスク等の記憶装置のことをいう。さらに「コンピュータ読み取り可能な記録媒体」とは、インターネット等のネットワークや電話回線等の通信回線を介してプログラムを送信する場合の通信線のように、短時間の間、動的にプログラムを保持するもの、その場合のサーバやクライアントとなるコンピュータシステム内部の揮発性メモリーのように、一定時間プログラムを保持しているものも含むものとする。また上記プログラムは、前述した機能の一部を実現するためのものであっても良く、さらに前述した機能をコンピュータシステムにすでに記録されているプログラムとの組み合わせで実現できるものであっても良い。

10

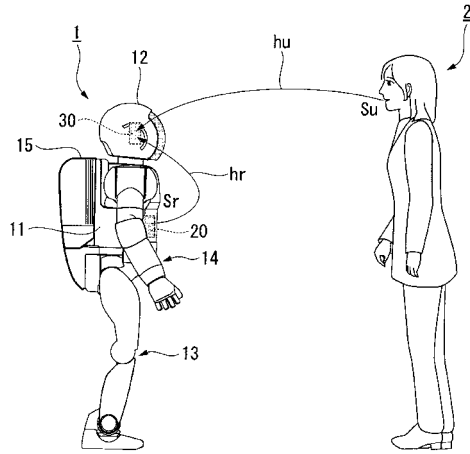
【符号の説明】

【0109】

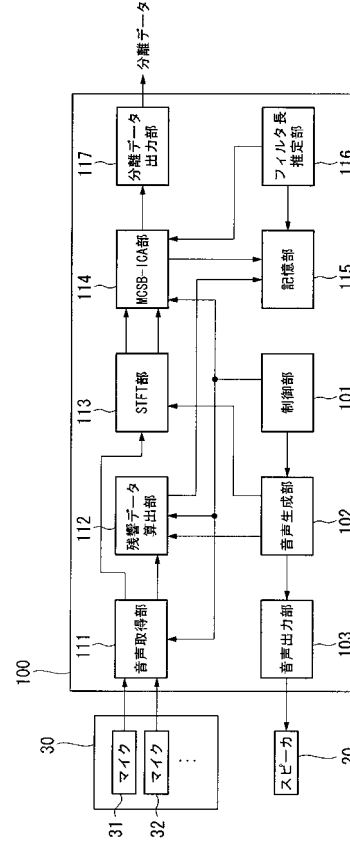
- 1・・・ロボット
- 20・・・スピーカ
- 30、31、32・・・マイク
- 100・・・残響抑圧装置
- 101・・・制御部
- 102・・・音声生成部
- 111・・・音声取得部
- 112・・・残響データ演算部
- 113・・・STFT部
- 114・・・MCMB-ICA部
- 115・・・記憶部
- 116・・・フィルタ長推定部
- 117・・・分離データ出力部
- 302・・・環境変化検出部

20

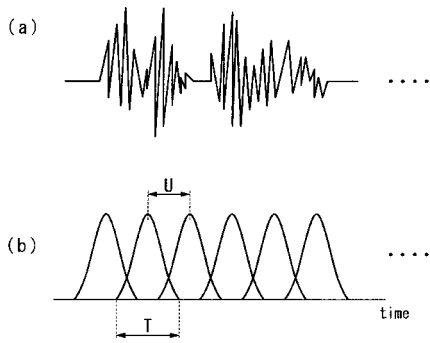
【図1】



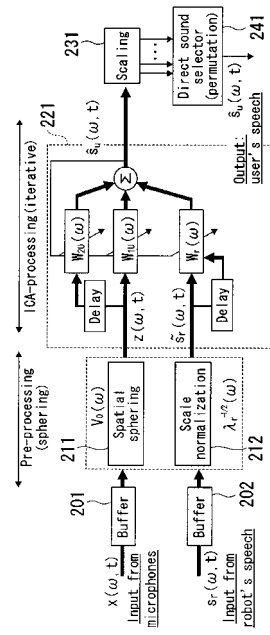
【図2】



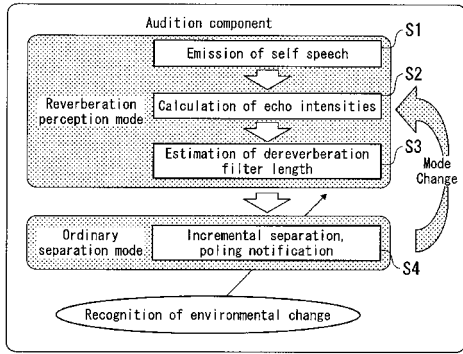
【図3】



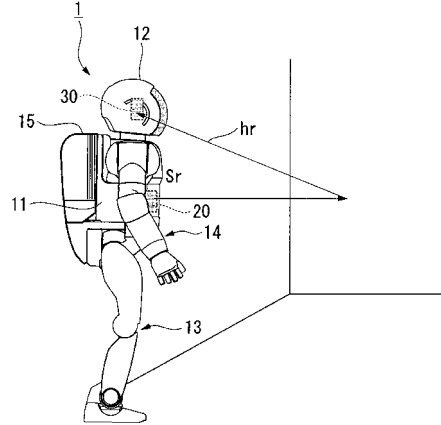
【図4】



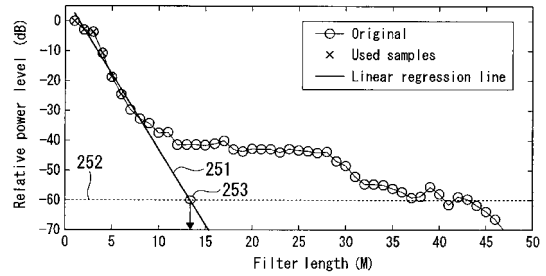
【 5 】



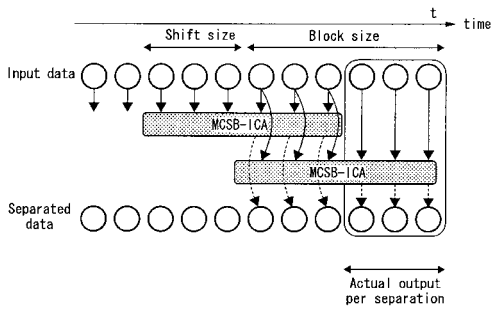
【 6 】



【 7 】



【 8 】



【 1 1 】

		Env. I (RT <sub>20</sub> 240 ms)			Env. II (RT <sub>20</sub> 670ms)		
		Mmax	20	30	50	30	40
w/o noise	Mean	14.0	13.7	13.2	35.0	35.3	35.4
	Std.	0.43	0.46	0.53	1.22	1.24	1.28
	with noise	Mean	14.2	14.0	13.6	36.1	36.3
	Std.	1.25	1.17	1.05	2.38	2.41	2.30

【 1 2 】

	Exp. B (non-barge-in)				Exp. C (barge-in)			
	no proc.	2s	2.5s	3s	no proc.	2s	2.5s	3s
Env. I (RT <sub>20</sub> 240ms)	74.3	76.9	78.5	78.2	28.2	67.8	70.2	71.7
Env. II (RT <sub>20</sub> 670ms)	26.1	63.9	66.8	69.2	11.0	37.1	41.2	43.3

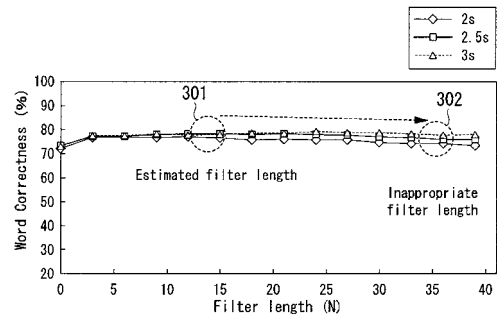
【 9 】

Impulse response	16-kHz sampling
Reverberation time (RT <sub>20</sub> )	240 and 670 ms
Distance and direction	1.5m and 0°, 45°, 90°, -45°, -90°
Number of microphones	Two (embedded in Robot's head)
SIFT analysis	Hanning: 32ms and shift: 12ms
Input wave data	[-1.0 1.0] normalized

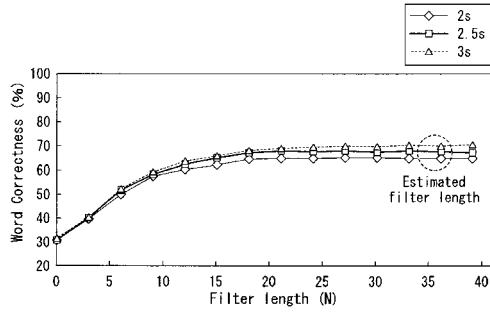
【 1 0 】

Test set	200 sentences
Training set	200 persons (150 sentences each)
Acoustic model	PTM-triphone: 3-state, HMM
Language model	Statistical, vocabulary size 20k
Speech analysis	Manning: 32ms and shift: 10ms
Features	MFCC 25 dim. (12+Δ12+ΔPow)

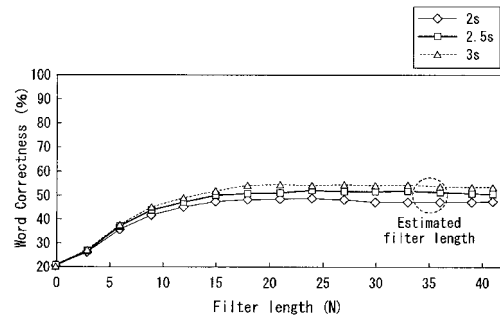
【 1 3 】



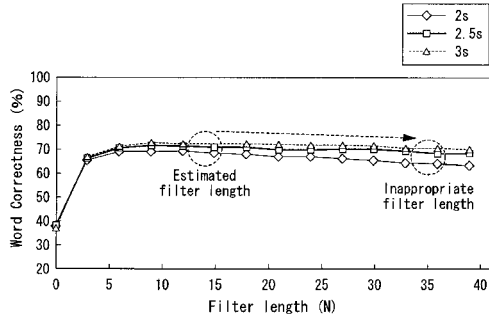
【 図 14 】



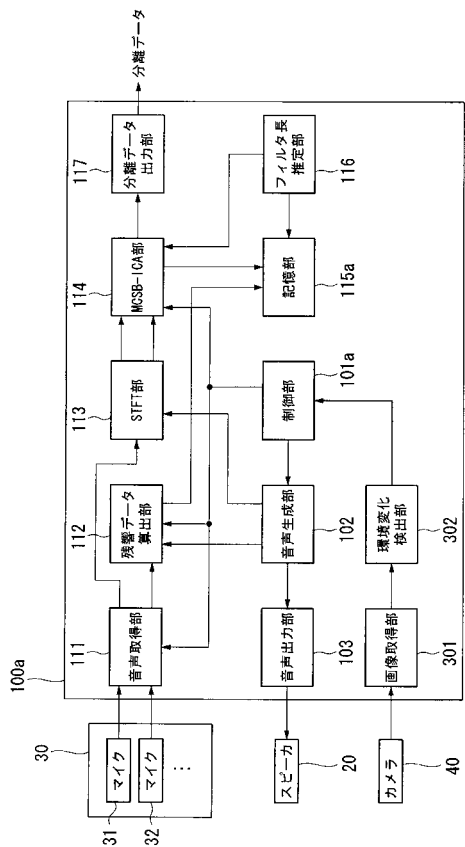
【 図 16 】



【 図 15 】



【 図 17 】



---

フロントページの続き

- (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
- (72)発明者 武田 龍  
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  
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

審査官 山下 剛史

- (56)参考文献 特開昭64 - 29094 (JP, A)  
特開2002 - 237770 (JP, A)  
特開平10 - 56406 (JP, A)  
特開昭64 - 29093 (JP, A)  
特開2009 - 276365 (JP, A)  
特開2009 - 159274 (JP, A)

(58)調査した分野(Int.Cl., DB名)

G10L 21/02 - 21/0308  
H04R 3/00 - 3/02