

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4825552号
(P4825552)

(45) 発行日 平成23年11月30日 (2011.11.30)

(24) 登録日 平成23年9月16日 (2011.9.16)

(51) Int. Cl.

F I

G 1 0 L 15/20 (2006.01)

G 1 0 L 15/20 3 7 0 D

G 1 0 L 21/02 (2006.01)

G 1 0 L 15/20 3 8 0

G 0 6 T 7/60 (2006.01)

G 1 0 L 21/02 1 0 1 B

G 1 0 L 21/02 1 0 2 B

G 1 0 L 21/02 2 0 3 Z

請求項の数 6 (全 17 頁) 最終頁に続く

(21) 出願番号 特願2006-68264 (P2006-68264)
 (22) 出願日 平成18年3月13日 (2006.3.13)
 (65) 公開番号 特開2007-248534 (P2007-248534A)
 (43) 公開日 平成19年9月27日 (2007.9.27)
 審査請求日 平成20年7月9日 (2008.7.9)

(73) 特許権者 504143441
 国立大学法人 奈良先端科学技術大学院大学
 奈良県生駒市高山町8916-5
 (73) 特許権者 000003207
 トヨタ自動車株式会社
 愛知県豊田市トヨタ町1番地
 (74) 代理人 110000110
 特許業務法人快友国際特許事務所
 (72) 発明者 猿渡 洋
 奈良県生駒市高山町8916-5 国立大
 学法人 奈良先端科学技術大学院大学内
 (72) 発明者 高谷 智哉
 奈良県生駒市高山町8916-5 国立大
 学法人 奈良先端科学技術大学院大学内
 最終頁に続く

(54) 【発明の名称】 音声認識装置、周波数スペクトル取得装置および音声認識方法

(57) 【特許請求の範囲】

【請求項 1】

対話者が話しかける音声から言語の内容を認識する音声認識装置であって、
 音を入力して、音信号に変換する複数の音入力手段と、
 音信号を周波数スペクトル（原周波数スペクトル）に変換する周波数変換手段と、
 原周波数スペクトルから、対話者の音声成分を含む1の周波数スペクトル（基本周波数
 スペクトル）を取得する音声成分検出手段と、
 複数の原周波数スペクトルから、フィルタ行列を用いる独立成分分析によって、雑音成
 分の1の周波数スペクトル（雑音周波数スペクトル）を取得する雑音成分推定手段と、
 基本周波数スペクトルから雑音周波数スペクトルを減算して、対話者の音声成分の周波
 数スペクトル（音声周波数スペクトル）を取得するスペクトル減算手段と、
 音声周波数スペクトルに基づいて、対話者が話しかけた言語の内容を認識する言語内容認
 識手段と、

対話者の方向を取得する方向取得手段と、

フィルタ行列の更新を行うか否かを判定する更新判定手段を備えており、

更新判定手段は、対話者の方向が変化した場合に、フィルタ行列の更新を行うと判定し

、

雑音成分推定手段は、更新判定手段によって更新を行うと判定された場合に、フィルタ
 行列の更新を行うことを特徴とする音声認識装置。

【請求項 2】

音声成分検出手段が、複数の原周波数スペクトルから、対話者の方向からの音声成分を強調した１の周波数スペクトルを基本周波数スペクトルとして取得することを特徴とする請求項１の音声認識装置。

【請求項３】

方向取得手段が、対話者を繰り返し撮影し、撮影された画像データを時刻と関連付ける撮像手段と、画像データから、対話者の方向を特定する方向特定手段を備えることを特徴とする請求項１または２の音声認識装置。

【請求項４】

更新判定手段が、対話者の方向がフィルタ行列の前の更新時から所定の角度以上変化した場合に、フィルタ行列の更新を行うと判定することを特徴とする請求項１から３の何れか一項の音声認識装置。

【請求項５】

対話者が話しかける音声の周波数スペクトルを取得する装置であって、
音を入力して、音信号に変換する複数の音入力手段と、
音信号を周波数スペクトル（原周波数スペクトル）に変換する周波数変換手段と、
原周波数スペクトルから、対話者の音声成分を含む１の周波数スペクトル（基本周波数スペクトル）を取得する音声成分検出手段と、

複数の原周波数スペクトルから、フィルタ行列を用いる独立成分分析によって、雑音成分の１の周波数スペクトル（雑音周波数スペクトル）を取得する雑音成分推定手段と、

基本周波数スペクトルから雑音周波数スペクトルを減算して、対話者の音声成分の周波数スペクトル（音声周波数スペクトル）を取得するスペクトル減算手段と、

対話者の方向を取得する方向取得手段と、

フィルタ行列の更新を行うか否かを判定する更新判定手段を備えており、

更新判定手段は、対話者の方向が変化した場合に、フィルタ行列の更新を行うと判定し

、雑音成分推定手段は、更新判定手段によって更新を行うと判定された場合に、フィルタ行列の更新を行うことを特徴とする周波数スペクトル取得装置。

【請求項６】

対話者が話しかける音声から言語の内容を認識する音声認識方法であって、
複数の音入力手段によって、入力される音を音信号に変換する音入力工程と、
音信号を周波数スペクトル（原周波数スペクトル）に変換する周波数変換工程と、
原周波数スペクトルから、対話者の音声成分を含む１の周波数スペクトル（基本周波数スペクトル）を取得する音声成分検出工程と、

複数の原周波数スペクトルから、フィルタ行列を用いる独立成分分析によって、雑音成分の１の周波数スペクトル（雑音周波数スペクトル）を取得する雑音成分推定工程と、

基本周波数スペクトルから雑音周波数スペクトルを減算して、対話者の音声成分の周波数スペクトル（音声周波数スペクトル）を取得するスペクトル減算工程と、

音声周波数スペクトルに基づいて、対話者が話しかけた言語の内容を認識する言語内容認識工程と、

対話者の方向を取得する方向取得工程と、

フィルタ行列の更新を行うか否かを判定する更新判定工程と、

更新判定工程において更新を行うと判定された場合に、フィルタ行列の更新を行う更新工程を備えており、

更新判定工程では、対話者の方向が変化した場合に、フィルタ行列の更新を行うと判定することを特徴とする音声認識方法。

【発明の詳細な説明】

【技術分野】

【０００１】

本発明は、対話者が話しかける音声を認識する装置と方法、および対話者が話しかける

10

20

30

40

50

音声の周波数スペクトルを取得する装置に関する。

【背景技術】

【0002】

人間が装置の動作を制御する際に、キーボードやレバーなどのインターフェースを操作することなく、音声を発することによって装置を制御する技術がある。このような技術においては、マイクロホンなどの音入力手段から入力される音声から、音声によって表現される言語の内容を認識し、認識された言語の内容に応じた制御が行われる。

【0003】

音入力手段から入力される音には、対話者が発した音声以外にも、周囲の雑音が混入する場合がある。周囲の雑音が混入した音に基づいて言語の内容の認識を行うと、誤認識を
10 起こし、装置の誤作動などを引き起こしてしまう。音声を認識する技術においては、雑音の影響をいかにして排除するかが重要である。

【0004】

入力される音から雑音成分を除去する技術が従来から開発されている。例えば特許文献1には、独立成分分析（ICA：Independent Component Analysis）を用いて、入力される音から雑音成分を除去する技術が開示されている。

【0005】

図5を参照しながら、ICAを用いた雑音成分の除去技術の概要を説明する。図5の音源502と音源504では、それぞれが別個に独立して音を発生する。マイクロホン506、508は、入力される音を音信号に変換する。マイクロホン506には、音源502
20 からの音と音源504からの音が重畳した音が入力される。マイクロホン508にも、音源502からの音と音源504からの音が重畳した音が入力される。図5では、音源502で発生する音を $s_1(t)$ と示し、音源504で発生する音を $s_2(t)$ と示す。また、マイクロホン506で取得される音を $x_1(t)$ と示し、マイクロホン508で取得される音を $x_2(t)$ と示す。ICAを用いる手法では、マイクロホン506と508で観測される観測信号 $x_1(t)$ と $x_2(t)$ を、フィルタ行列510を用いて、出力信号 $y_1(t)$ と $y_2(t)$ に分離する。フィルタ行列510は、出力信号 $y_1(t)$ 、 $y_2(t)$ が統計的に独立となるように、フィルタ最適化手段512によって、その係数が最適化される。フィルタ行列510の最適化には、例えば出力信号 $y_1(t)$ 、 $y_2(t)$ についてのコスト関数を最小化する手法などが知られている。
30

音源502で発生する音 $s_1(t)$ と音源504で発生する音 $s_2(t)$ は統計的に独立であるから、 $x_1(t)$ と $x_2(t)$ から抽出された統計的に独立な出力信号 $y_1(t)$ と $y_2(t)$ は、一方が音源502で発生した音 $s_1(t)$ であり、他方が音源504で発生した音 $s_2(t)$ であると推定される。

【0006】

上記では2つの音源502、504で発生した音を分離する例を説明しているが、同様の手法によって、3つ以上の音源からの音が重畳した音が入力される場合に、認識の対象としたい特定の1つの音源からの音と、それ以外の音源からの音が重畳した音に分離することができる。これによって、音声認識の対象としたい音源からの音（対話者の音声）と、それ以外の音源からの音が重畳した音（雑音）を分離することができる。
40

【0007】

また、特許文献2には、スペクトル・サブトラクション（SS：Spectral Subtraction）法を用いて、入力される音から雑音成分を除去する技術が開示されている。

【0008】

図6を参照しながら、SS法を用いた雑音成分の除去技術の概要を説明する。SS法では、同一の面に沿って所定の間隔 d で配置された複数のマイクロホン604、606、
50 ・ ・ ・を用いる。間隔 d に比べて十分に離れた位置にある音源602で発生した音は、マイクロホン604、606、
・ ・ ・の近傍ではほぼ平面波として伝播して、マイクロホン604、606、
・ ・ ・に到達する。マイクロホン604に到達する音は経路610に沿って伝播してきており、マイクロホン606に到達する音は経路612に沿って伝播してき

ている。平面波として伝播してくる音は、波面 6 1 6 や波面 6 1 8 において、同一の位相を備えている。従って、音源 6 0 2 がマイクロホン 6 0 4、6 0 6、・・・から見て角度 θ の方向にある場合、経路 6 1 0 に沿ってある時点でマイクロホン 6 0 4 に到達した音は、経路 6 1 2 に沿って点 6 2 0 まで到達しており、その後さらに $d \cos \theta$ の長さの経路を伝播してから、マイクロホン 6 0 6 の点 6 2 2 に到達する。従って、隣接するマイクロホン 6 0 4 と 6 0 6 では、音の伝播速度を c とすると、到来時間差 $t = d \cos \theta / c$ で、それぞれ音が到達する。

【0009】

上記のような音源の方向と音の到来時間差についての関係を利用して、複数のマイクロホン 6 0 4、6 0 6、・・・に入力される音に基いて、特定の方向からの音を強調したり、逆に特定の方向からの音を抑圧したりすることができる。

10

例えば図 6 に示す例で、 θ 方向からの音を強調したい場合には、マイクロホン 6 0 6 で取得された音信号 6 2 6 と、マイクロホン 6 0 4 で取得された音信号 6 2 4 を $t = d \cos \theta / c$ だけ遅延させた信号の和を算出することで、方向 θ からの音を強調した音信号を得ることができる。

また、マイクロホン 6 0 6 で取得された音信号 6 2 6 と、マイクロホン 6 0 4 で取得された音信号 6 2 4 を $t = d \cos \theta / c$ だけ遅延させて正負を反転させた信号の和を算出し、この信号にスペクトル形状補正フィルタを適用することで、方向 θ からの音を抑圧した音信号を得ることができる。

上記のようにして、マイクロホン 6 0 4、6 0 6、・・・の配置と音の方向に応じた遅延時間をそれぞれ設定しておく。そして、それぞれのマイクロホン 6 0 4、6 0 6、・・・から入力される音信号を、音の方向に応じた遅延時間だけ遅延させてから和を算出する遅延和アレーを用いることで、特定の方向からの音を強調した信号を取得したり、逆に特定の方向からの音を抑圧した信号を取得したりすることができる。図 7 に、上記のような遅延和アレーを用いた場合の、特定の方向からの音を強調する指向特性 7 0 2 と、特定の方向からの音を抑圧する指向特性 7 0 4 の例を示す。図 7 では、 $\theta = 90^\circ$ の方向、すなわちマイクロホン 6 0 4、6 0 6、・・・の正面の方向からの音を強調したり抑圧したりする場合の指向特性を示している。

20

【0010】

上記において、対話者の方向からの音を抑圧した信号は、対話者の方向以外からの音を検出しており、雑音成分と推定することができる。SS 法では、対話者の方向からの音を強調した周波数スペクトルと、雑音成分の周波数スペクトルをそれぞれ特定し、両者の差をとることによって、雑音成分が除去された音声の周波数スペクトルを取得する。雑音成分が除去された周波数スペクトルから、音声認識の対象としたい音源からの音（対話者の音声）の特徴を把握することができる。

30

【0011】

【特許文献 1】特開 2 0 0 4 - 0 6 9 7 7 2 号公報

【特許文献 2】特開 2 0 0 1 - 1 0 0 8 0 0 号公報

【発明の開示】

【発明が解決しようとする課題】

40

【0012】

上述のような I C A を用いた雑音成分の除去技術では、対話者の位置や、それぞれのマイクロホンの位置や、それぞれのマイクロホンの感度特性などが未知であっても、実際の状況に応じてフィルタ行列が最適化されるため、入力される音から雑音成分を除去することができる。従って、対話者が移動して位置が変化したり、マイクロホンの取付け位置や感度特性にばらつきがあったりしても、それに合わせてフィルタ行列が最適化されるため、入力される音から雑音成分を除去することができる。ロバストな音声認識システムを構築することができる。

【0013】

しかしながら、I C A による雑音成分の除去技術には、演算負荷が高いという問題があ

50

る。ICAを用いる場合、最終的な出力信号が統計的に独立となるように、常にフィルタ行列を最適化する必要がある。例えば音源の位置が変わった場合には、フィルタ行列は新たに最適化しなければならない。従って、ICAによる雑音成分の除去では、音源の移動の有無に関わらず、フィルタ行列を所定の時間間隔で繰り返し更新して、フィルタ行列が最適化された状態を維持する必要がある。このようにフィルタ行列の更新を常時行っていると、計算の負荷が増大して、処理に要する時間が長いものになってしまう。処理に要する時間が長いと、対話者への応答が遅れ、対話者は不快感を覚えてしまう。

【0014】

上述したSS法を用いた雑音成分の除去技術では、フィルタ行列の最適化のような複雑な計算を行う必要がないため、演算の負荷はそれほど高いものではない。従って、対話者が話しかけてから、短時間で雑音成分を除去することができる。さらにSS法では、雑音成分を除去するのみではなく、対話者の方向からの音を強調しているため、ICAを用いる場合に比べて、対話者の音声の特徴をより鮮明に抽出することができる。

【0015】

しかしながら、SS法による雑音成分の除去技術は、対話者の位置や、マイクロホンの位置や、マイクロホンの感度特性の変動に影響を受けるという問題がある。特に、マイクロホンの感度特性のばらつきが、雑音成分の除去に大きな影響を及ぼす。

図8と図9は、マイクロホンの感度特性のばらつきが、遅延和アレーによって実現される指向特性に及ぼす影響を示している。図8はマイクロホンの感度特性のばらつきが、特定の方向からの音を強調する処理の際に用いられる指向特性に及ぼす影響を示す。分布702は全てのマイクロホンが同じ感度特性を持つ場合の理想的な指向特性を示す。分布810と分布812は、マイクロホンに感度特性のばらつきがある場合の指向特性を示す。ここでは一例として、マイクロホンの感度特性に ± 2 dBのばらつきがある場合の指向特性を分布810で示し、マイクロホンの感度特性に ± 4 dBのばらつきがある場合の指向特性を分布812で示す。図8から明らかなように、音声を強調する処理の際には、マイクロホンの感度特性のばらつきによって、わずかに指向特性が鈍化するものの、大きな影響はない。一方、図9はマイクロホンの感度特性のばらつきが、特定の方向からの音を抑圧する処理の際に用いられる指向特性に及ぼす影響を示す。分布704は全てのマイクロホンが同じ感度特性を持つ場合の理想的な指向特性を示す。分布910と分布912は、マイクロホンに感度特性のばらつきがある場合の指向特性を示す。ここでは一例として、マイクロホンの感度特性に ± 2 dBのばらつきがある場合の指向特性を分布910で示し、マイクロホンの感度特性に ± 4 dBのばらつきがある場合の指向特性を分布912で示す。図9から明らかなように、特定の方向からの音を抑圧する処理の際に用いられる指向特性は、マイクロホンの感度特性のばらつきによって、大きな影響を受ける。マイクロホン間の感度特性にばらつきがあると、特定の方向からの音について、ほとんど抑圧することができなくなってしまう。

【0016】

上記のように、遅延和アレーによって特定の方向からの音声を抑圧する場合、マイクロホンの感度特性のばらつきが大きな影響を及ぼす。マイクロホンの感度特性にばらつきがあると、音声認識の対象としたい音声まで雑音成分に含ませてしまうことになる。従って、対話者の方向からの音を強調したスペクトルから雑音成分のスペクトルを減算する際に、本来は減算すべきでない対話者の方向からの音の成分についてまで減らすことになってしまう。正確に雑音成分の除去を行うことが困難となる。

【0017】

上述のように、ICAを用いる技術と、SS法を用いる技術には、それぞれ一長一短がある。ICAを用いる場合には、マイクロホンの感度特性のばらつきは何ら影響しないが、計算負荷が高く、処理が遅くなる。SS法を用いる場合には、計算負荷が軽いものの、マイクロホンの感度特性のばらつきによって、正確な雑音成分の除去が困難になってしまう。処理の負荷が軽く、なおかつマイクロホンの感度特性のばらつきの影響を受けずに、対話者の音声を鮮明に抽出することが可能な技術が待望されている。

【 0 0 1 8 】

本発明では上記課題を解決する。本発明では、少ない計算負荷で、マイクロホンの感度特性のばらつきの影響を受けずに、対話者の音声を鮮明に抽出して、正確な音声認識を行うことが可能な技術を提供する。

【課題を解決するための手段】

【 0 0 1 9 】

本発明は装置として具現化される。本発明の装置は、対話者が話しかける音声から言語の内容を認識する音声認識装置である。その装置は、音を入力して音信号に変換する複数の音入力手段と、音信号を周波数スペクトル（原周波数スペクトル）に変換する周波数変換手段と、原周波数スペクトルから対話者の音声成分を含む１の周波数スペクトル（基本周波数スペクトル）を取得する音声成分検出手段と、複数の原周波数スペクトルからフィルタ行列を用いる独立成分分析によって雑音成分の１の周波数スペクトル（雑音周波数スペクトル）を取得する雑音成分推定手段と、基本周波数スペクトルから雑音周波数スペクトルを減算して対話者の音声成分の周波数スペクトル（音声周波数スペクトル）を取得するスペクトル減算手段と、音声周波数スペクトルに基いて対話者が話しかけた言語の内容を認識する言語内容認識手段を備えている。

10

【 0 0 2 0 】

上記の音声認識装置は、複数の音入力手段からの音信号に基いて、対話者の音声成分を含む周波数スペクトル（基本周波数スペクトル）と、雑音成分の周波数スペクトル（雑音周波数スペクトル）をそれぞれ取得して、両者の差を取ることによって、対話者の音声成分の周波数スペクトル（音声周波数スペクトル）を取得する。基本周波数スペクトルは、例えば複数の原周波数スペクトルのうちの１つであってもよいし、遅延和アレーによって対話者の方向からの音声を強調した周波数スペクトルであってもよい。雑音周波数スペクトルは、フィルタ行列を用いるICAによって、取得することができる。上記のようにして得られる音声周波数スペクトルは、対話者の音声成分を含んでいる基本周波数スペクトルから、雑音周波数スペクトルを除去したものであるから、対話者の音声を鮮明に抽出したものである。上記の音声認識装置は、スペクトル減算手段で取得された音声周波数スペクトルに基いて、言語内容の認識を行う。このような構成とすることによって、処理に要する時間が短く、かつ誤認識を起こしにくい音声認識装置を実現することができる。

20

【 0 0 2 1 】

上記の音声認識装置によれば、複数の音声入力手段の間で感度特性にばらつきがあっても、従来のSS法を用いた技術とは異なり、雑音成分推定手段は正確に雑音成分を推定することができる。これによって、対話者の音声成分の周波数スペクトルを鮮明に抽出して、言語内容の認識を行うことができる。言語内容の認識の精度を向上することができる。

30

【 0 0 2 2 】

上記の音声認識装置は、フィルタ行列の更新を行うか否かを判定する更新判定手段をさらに備えており、雑音成分推定手段が、更新判定手段によって更新を行うと判定された場合に独立成分分析で用いるフィルタ行列の更新を行うことが好ましい。

【 0 0 2 3 】

上記の音声認識装置によれば、雑音成分推定手段における独立成分分析で用いるフィルタ行列が必要な時にのみ更新されるため、雑音成分推定手段での演算の負荷が軽減される。音声認識に係る処理時間を短縮することができる。

40

【 0 0 2 4 】

上記の音声認識装置は、対話者の方向を取得する方向取得手段をさらに備えており、音声成分検出手段が、複数の原周波数スペクトルから対話者の方向からの音声成分を強調した１の周波数スペクトルを基本周波数スペクトルとして取得することがさらに好ましい。

【 0 0 2 5 】

対話者の方向を取得することができれば、複数の原周波数スペクトルによる遅延和アレーを用いることで、対話者の方向からの音声成分を強調した周波数スペクトルを得ることができる。対話者の方向の取得は、例えばセンサ等を用いて対話者の位置を検出してもよ

50

いし、予め対話者の位置を制限しておいて、その位置を記憶しておいてもよい。

上記の音声認識装置によれば、対話者の方向からの音声成分を強調した周波数スペクトルが基本周波数スペクトルとして取得されるため、対話者の音声成分をより鮮明に抽出することができる。言語内容の認識の精度をさらに向上することができる。

【0026】

上記の音声認識装置においては、方向取得手段が、対話者を繰り返し撮影し撮影された画像データを時刻と関連付ける撮像手段と、画像データから対話者の方向を特定する方向特定手段を備えており、更新判定手段が、対話者の方向が前回の更新時から所定の角度以上変化した場合にフィルタ行列の更新を行うと判定することがさらに好ましい。

【0027】

10

対話者の方向は、例えば複数の撮像手段からの画像データを用いてステレオ視の原理によって計算することができる。

上記の音声認識装置によれば、画像データに基いて対話者の正確な方向を取得することができるため、音声成分強調手段において対話者の方向からの音声成分をよりの確に強調することができる。さらに、上記の音声認識装置によれば、対話者の方向が所定の角度以上変化した時点で雑音成分推定手段のフィルタ行列の更新を行うため、不要な更新処理を行うことがなく、かつ必要な更新処理は確実に行うため、処理の負荷をさらに軽減して対話者の音声成分をより鮮明に抽出することができる。言語内容の認識の精度をさらに向上することができる。

【0028】

20

上記したフィルタ行列の更新は、音声認識装置と対話者との位置関係が変化した場合のみ行うのではなく、音声認識装置と対話者の周囲の環境が変化した場合に行うことで、独立成分分析に用いるフィルタ行列を適切に更新することができる。

すなわち、上記の音声認識装置は、対話者を繰り返し撮影し撮影された画像データを時刻と関連付ける撮像手段と、画像データから周囲の環境の変化を認識する環境認識手段をさらに備えており、更新判定手段が、周囲の環境が変化した場合にフィルタ行列の更新を行うと判定することも好ましい。

【0029】

周囲の環境の変化は、例えば画像データから対話者の輪郭を抽出し、対話者の輪郭に基づいて対話者の映像と周囲の環境の映像を識別し、周囲の環境の映像の経時的変化から認識することができる。

30

上記の音声認識装置によれば、雑音成分推定手段における独立成分分析で用いるフィルタ行列が必要な時にのみ更新されるため、雑音成分推定手段での演算の負荷が軽減される。音声認識に係る処理時間を短縮することができる。

【0030】

本発明は、対話者の音声の周波数スペクトルを取得する装置としても具現化される。本発明の他の1つの装置は、音を入力して音信号に変換する複数の音入力手段と、音信号を周波数スペクトル（原周波数スペクトル）に変換する周波数変換手段と、原周波数スペクトルから対話者の音声成分を含む1の周波数スペクトル（基本周波数スペクトル）を取得する音声成分検出手段と、複数の原周波数スペクトルからフィルタ行列を用いる独立成分分析によって雑音成分の1の周波数スペクトル（雑音周波数スペクトル）を取得する雑音成分推定手段と、基本周波数スペクトルから雑音周波数スペクトルを減算して対話者の音声成分の周波数スペクトル（音声周波数スペクトル）を取得するスペクトル減算手段を備えている。

40

【0031】

本発明は方法としても具現化される。本発明の方法は、対話者が話しかける音声から言語の内容を認識する音声認識方法である。その方法は、複数の音入力手段によって入力される音を音信号に変換する工程と、音信号を周波数スペクトル（原周波数スペクトル）に変換する工程と、原周波数スペクトルから対話者の音声成分を含む1の周波数スペクトル（基本周波数スペクトル）を取得する工程と、複数の原周波数スペクトルからフィルタ行

50

列を用いる独立成分分析によって雑音成分の１の周波数スペクトル（雑音周波数スペクトル）を取得する工程と、基本周波数スペクトルから雑音周波数スペクトルを減算して対話者の音声成分の周波数スペクトル（音声周波数スペクトル）を取得する工程と、音声周波数スペクトルに基いて対話者が話しかけた言語の内容を認識する工程を備えている。

【発明の効果】

【００３２】

本発明によれば、少ない計算負荷で、マイクロホンの感度特性のばらつきの影響を受けずに、対話者の音声を鮮明に抽出して、正確な音声認識を行うことができる。

【発明を実施するための最良の形態】

【００３３】

以下に発明を実施するための最良の形態を列記する。

（形態１）音声成分検出手段は、複数の原周波数スペクトルによる遅延和アレーを用いて、対話者の方向からの音声成分を強調した１の周波数スペクトルを基本周波数スペクトルとして取得する。

【実施例】

【００３４】

本実施例では、図１に例示する音声認識装置１００において、対話者Ｖが話しかける音声を認識する例を説明する。音声認識装置１００は、例えばショールームやイベント会場に配置された案内ロボットであり、案内を求めて話しかけてくる来場者（対話者）Ｖが話しかける音声を認識する。

【００３５】

音声認識装置１００は、頭部１０２の前方に並んで配置された右カメラ１０４と左カメラ１０６と、胴体部１０８の前方の集音部１１０に所定の間隔で並んで配置されたマイクロホン１１２ａ、１１２ｂ、１１２ｃ、１１２ｄ、１１２ｅおよび１１２ｆと、頭部１０２の前方に配置されたスピーカ１１６と、右カメラ１０４、左カメラ１０６、マイクロホン１１２ａ、１１２ｂ、１１２ｃ、１１２ｄ、１１２ｅおよび１１２ｆ、スピーカ１１６と通信可能なコントローラ１１４を備えている。

【００３６】

右カメラ１０４と左カメラ１０６は、一般的なＣＣＤカメラである。右カメラ１０４と左カメラ１０６は、所定の時間間隔で同時に撮影を実施し、撮影された画像データを撮影時刻と関連付けてコントローラ１１４へ出力する。

【００３７】

マイクロホン１１２ａ、１１２ｂ、１１２ｃ、１１２ｄ、１１２ｅおよび１１２ｆは、入力される音声によって振動板に加えられる音圧を検知し、検知した音圧に応じた電圧値を内蔵されたアンプによって増幅し、コントローラ１１４へ出力する。

【００３８】

スピーカ１１６は、コントローラ１１４から送信される信号をアンプによって増幅し、増幅された電流の変動に応じて振動板を振動させ、音声を出力する。

【００３９】

図２はコントローラ１１４の構成を示すブロック図である。コントローラ１１４は、処理装置（ＣＰＵ）、記憶装置（光学記憶媒体、磁気記憶媒体、あるいはＲＡＭやＲＯＭといった半導体メモリ等）、入出力装置、演算装置などから構成されているコンピュータ装置である。コントローラ１１４は、右カメラ１０４、左カメラ１０６から入力される画像データと、マイクロホン１１２ａ、１１２ｂ、１１２ｃ、１１２ｄ、１１２ｅおよび１１２ｆから入力される音信号に基いて、対話者Ｖが話しかける言語の内容を認識して、その内容に対する返答をスピーカ１１６から音声で出力する。コントローラ１１４は機能的に、画像認識部２０２、更新判定部２０４、Ａ／Ｄ変換部２０６、周波数変換部２０８、音声成分強調部２１０、雑音成分推定部２１２、メル周波数変換部２１４および２１６、スペクトル減算部２１８、特徴量計算部２２０、言語認識部２２２、応答制御部２２４、Ｄ／Ａ変換部２２６を備えている。

10

20

30

40

50

【 0 0 4 0 】

画像認識部 2 0 2 は、右カメラ 1 0 4 と左カメラ 1 0 6 から出力される画像データに基づいて、対話者 V の位置を認識する。対話者 V の位置は、右カメラ 1 0 4 と左カメラ 1 0 6 のそれぞれの画像データにおいて対話者 V の輪郭を抽出し、輪郭を抽出された対話者 V と音声認識装置 1 0 0 との相対的な位置関係をステレオ視の原理によって算出することで、算出することができる。画像認識部 2 0 2 は、対話者 V の位置から対話者 V の方向を特定する。画像認識部 2 0 2 は、特定された対話者 V の方向を、時刻と関連付けて、更新判定部 2 0 4 と音声成分強調部 2 1 0 に出力する。

【 0 0 4 1 】

更新判定部 2 0 4 は、対話者 V の方向に基いて、音声成分強調部 2 1 0 において音を強調する方向の更新と、雑音成分推定部 2 1 2 において雑音成分の推定に用いるフィルタ行列の更新を行うか否かを判定する。

更新判定部 2 0 4 は、前回更新を行った際の対話者 V の方向を保持している。更新判定部 2 0 4 は、画像認識部 2 0 2 から新たに対話者 V の方向が入力されると、前回更新を行った際の対話者 V の方向との比較を行い、所定の角度以上変化しているか否かを評価する。対話者 V の方向が、前回の更新時点から所定の角度以上変化している場合に、更新判定部 2 0 4 は更新が必要であると判定する。更新が必要であると判定した場合、更新判定部 2 0 4 は、音声成分強調部 2 1 0 と雑音成分推定部 2 1 2 に、更新指示を出力する。

【 0 0 4 2 】

A / D 変換部 2 0 6 は、マイクロホン 1 1 2 a、1 1 2 b、1 1 2 c、1 1 2 d、1 1 2 e および 1 1 2 f から入力されるそれぞれの音信号を、A / D 変換してデジタル化する。図 2 では図示の簡略化のために A / D 変換部 2 0 6 は 1 つのブロックとして図示されているが、本実施例の音声認識装置 1 0 0 は、マイクロホン 1 1 2 a、1 1 2 b、1 1 2 c、1 1 2 d、1 1 2 e および 1 1 2 f のそれぞれに対応する A / D 変換手段を、並列に処理可能な状態で備えている。A / D 変換部 2 0 6 は、デジタル化された音信号を音データとして周波数変換部 2 0 8 へ出力する。

【 0 0 4 3 】

周波数変換部 2 0 8 は、A / D 変換部 2 0 6 から入力される音データのそれぞれについて、周波数スペクトルの時系列を特定する。周波数変換部 2 0 8 は、まず音データのフレーム化処理を行い、次いで各フレームについての高速フーリエ変換を行って、周波数スペクトルの時系列を特定する。図 1 0 に音データのフレーム化処理と、各フレームの音データの周波数スペクトルの時系列を特定する様子を示す。本実施例では、フレームの長さは 2 0 m s であり、フレーム間隔は 1 0 m s である。図 1 0 に示すように、音データ 1 0 0 2 についてフレーム F 1、F 2、F 3、・・・が規定される。周波数変換部 2 0 8 は、フレーム F 1、F 2、F 3、・・・のそれぞれにおける音データ 1 0 0 2 の周波数スペクトル f_1 、 f_2 、 f_3 、・・・を特定する。周波数スペクトルは、周波数に対する振幅の分布として与えられる。周波数スペクトルの特定は、例えば高速フーリエ変換を用いて行うことができる。

【 0 0 4 4 】

なお図 2 では図示の簡略化のために周波数変換部 2 0 8 は 1 つのブロックとして図示されているが、本実施例の音声認識装置 1 0 0 は、マイクロホン 1 1 2 a、1 1 2 b、1 1 2 c、1 1 2 d、1 1 2 e および 1 1 2 f の入力信号をデジタル化した音データのそれぞれに対応する周波数変換手段を、並列に処理可能な状態で備えている。周波数変換部 2 0 8 は、特定された周波数スペクトルの時系列を音声成分強調部 2 1 0 と雑音成分推定部 2 1 2 に出力する。

【 0 0 4 5 】

音声成分強調部 2 1 0 は、マイクロホン 1 1 2 a、1 1 2 b、1 1 2 c、1 1 2 d、1 1 2 e および 1 1 2 f のそれぞれに対応する周波数スペクトルの時系列から、対話者 V の方向から到来した音を強調した周波数スペクトル（以下では基本周波数スペクトルと呼ぶ）の時系列を算出する。音声成分強調部 2 1 0 は、音を強調する方向を保持している。音

10

20

30

40

50

声成分強調部 210 は、更新判定部 204 から更新指示が送信される度に、音を強調する方向を、画像認識部 202 から入力される対話者 V の方向に更新する。

【0046】

音声成分強調部 210 は、マイクロホン 112 a、112 b、112 c、112 d、112 e および 112 f のそれぞれに対応する周波数スペクトルの時系列から、遅延和アレーを用いて、基本周波数スペクトルを算出する。具体的には、音を強調する方向から、各マイクロホンに対応する遅延時間を特定し、特定された遅延時間に相当するフレーム数を特定する。音声成分強調部 210 は、マイクロホン 112 a、112 b、112 c、112 d、112 e および 112 f のそれぞれに対応する周波数スペクトルについて、それぞれについて特定されたフレーム数だけオフセットさせて、和を算出する。音声成分強調部 210 は、算出された基本周波数スペクトルの時系列を、メル周波数変換部 214 へ出力する。

10

【0047】

雑音成分推定部 212 は、マイクロホン 112 a、112 b、112 c、112 d、112 e および 112 f のそれぞれに対応する周波数スペクトルの時系列から、対話者 V の音声以外の音の周波数スペクトル（雑音周波数スペクトル）の時系列を算出する。

雑音成分推定部 212 は、マイクロホン 112 a、112 b、112 c、112 d、112 e および 112 f のそれぞれに対応する周波数スペクトルから、統計的に独立な 2 の周波数スペクトルを算出する、フィルタ行列を保持している。雑音成分推定部 212 は、統計的に独立な 2 の周波数スペクトルが算出されると、両者のうちの一方を雑音周波数スペクトルとして選択する。雑音周波数スペクトルの選択は、例えば雑音の周波数スペクトルとして典型的な周波数スペクトル形状を予め記憶しておき、その典型的な周波数スペクトル形状との類似性が高い方を、雑音周波数スペクトルとして選択することができる。雑音成分推定部 212 は、推定された雑音周波数スペクトルの時系列を、メル周波数変換部 216 へ出力する。

20

なお、雑音成分推定部 212 は、更新判定部 204 から更新指示が入力されると、保持されているフィルタ行列の最適化処理を行う。フィルタ行列の最適化処理は、例えばフィルタ行列によって分離される 2 の周波数スペクトルについてのコスト関数を最小化する手法によって行うことができる。

【0048】

30

メル周波数変換部 214 は、音声成分強調部 210 から出力される強調周波数スペクトルの時系列を、メル周波数に関するスペクトルの時系列へ変換する。メル周波数は音の高低に対する人間の感覚を示す尺度であって、周波数 f からメル周波数 $Mel(f)$ への変換は、次の関係式を用いて行うことができる。

【0049】

【数 1】

$$Mel(f) = 2595 \log_{10} \left(1 + \frac{f}{700} \right)$$

【0050】

40

メル周波数変換部 214 は、基本周波数スペクトルをメル周波数に変換したスペクトル（基本メル周波数スペクトル）の時系列を、スペクトル減算部 218 へ出力する。

【0051】

メル周波数変換部 216 は、雑音成分推定部 212 から出力される雑音周波数スペクトルの時系列を、メル周波数に関するスペクトルの時系列へ変換する。メル周波数変換部 216 は、雑音周波数スペクトルをメル周波数に変換したスペクトル（雑音メル周波数スペクトル）の時系列を、スペクトル減算部 218 へ出力する。

【0052】

スペクトル減算部 218 は、メル周波数変換部 214 から出力された基本メル周波数スペクトルの時系列から、メル周波数変換部 216 から出力された雑音メル周波数スペクトル

50

ルの時系列を減じて、雑音成分が除去された対話者Vの音声のメル周波数スペクトル（音声メル周波数スペクトル）の時系列を算出する。スペクトル減算部218は、特定された音声メル周波数スペクトルの時系列を特徴量計算部220へ出力する。

【0053】

特徴量計算部220は、スペクトル減算部218から出力される音声メル周波数スペクトルの時系列から、対話者Vが話しかけた音声の特徴量の時系列を算出する。本実施例では、特徴量として、メル周波数ケプストラム係数（MFCC）を扱う。特徴量計算部220は、計算されたMFCCの時系列を言語認識部222へ出力する。

【0054】

言語認識部222は、特徴量計算部220から出力されるMFCCの時系列から、対話者Vが話しかけた言語の内容を認識する。言語認識部222は、MFCCと音素との対応を示す対応表を保持しており、入力されるMFCCの時系列から、対話者Vが発した音素の時系列を特定する。さらに言語認識部222は、特定された音素の時系列から、対話者Vが話しかけた言語の内容を示す文字列を特定する。言語認識部222は、特定された文字列を応答制御部224へ出力する。

10

【0055】

応答制御部224は、言語認識部222で特定された文字列から、対話者Vが話しかけた言語の内容に対する適切な返答を生成する。応答制御部224は、対話者Vが話しかける言語の内容を示す文字列と、その内容に対する適切な返答を表現する音声データの対応を示す対応表を保持しており、入力される文字列をキーとして、適切な返答の音声データを検索する。応答制御部224は、検索された適切な返答の音声データをD/A変換部226へ出力する。

20

【0056】

D/A変換部226は、応答制御部224から入力される音声データをD/A変換して、スピーカ116へ出力する。これによって、対話者Vが話しかけた言語の内容に応じた適切な返答がスピーカ116から音声で出力される。

【0057】

図3と図4に示すフローチャートを用いて、音声認識装置100の動作を説明する。音声認識装置100は、図3のフローチャートに示す音声認識処理を常に行っている。それと並行して、音声認識装置100は、所定の時間間隔で、図4のフローチャートに示す更新判定処理を実行する。

30

【0058】

図3の音声認識処理を開始すると、ステップS302において、マイクロホン112a、112b、112c、112d、112eおよび112fから音信号が入力され、それらの音信号をA/D変換部206で音データに変換する。

ステップS304では、周波数変換部208が、音データの周波数スペクトルを特定する。

ステップS306では、音声成分強調部210が、基本周波数スペクトルを算出する。

ステップS308では、雑音成分推定部212が、雑音周波数スペクトルを算出する。

なお本実施例の音声認識装置100では、ステップS306とステップS308は、ステップS304の後に、同時に並行して実施される。ステップS306とステップS308の両方の処理が終了した後、処理はステップS310へ移行する。

40

ステップS310では、基本周波数スペクトルと、雑音周波数スペクトルを、それぞれメル周波数に関するスペクトルに変換する。

ステップS312では、スペクトル減算部218が、基本メル周波数スペクトルから、雑音メル周波数スペクトルを減算（スペクトル・サブトラクション）して、音声メル周波数スペクトルを算出する。

ステップS314では、特徴量計算部220が、音声メル周波数スペクトルから、対話者Vの音声の特徴量であるMFCCを計算する。

ステップS316では、言語認識部222が、MFCCから言語の内容を認識して対応

50

する文字列を特定する。

ステップS 3 1 8では、応答制御部2 2 4が、特定された文字列に応じた適切な返答を生成する。

ステップS 3 2 0では、スピーカ1 1 6から、返答が音声で出力される。

ステップS 3 2 0の後、処理はステップS 3 0 3へ移行し、上述の処理を繰り返し実行する。

【0 0 5 9】

図4の更新判定処理が開始されると、ステップS 4 0 2において、右カメラ1 0 4と左カメラ1 0 6から画像データを取得する。

ステップS 4 0 4では、画像認識部2 0 2が、取得された画像データから対話者Vの方向を特定する。

ステップS 4 0 6では、更新判定部2 0 4が、対話者Vの方向から、音声成分強調部2 1 0において音を強調する方向と、雑音成分推定部2 1 2におけるフィルタ行列の、更新が必要か否かを判断する。更新が必要と判定された場合（ステップS 4 0 6でYESの場合）、処理はステップS 4 0 8へ進む。更新が不要と判定された場合（ステップS 4 0 6でNOの場合）、音声成分強調部2 1 0において音を強調する方向と、雑音成分推定部2 1 2のフィルタ行列について、更新をすることなく、図4の更新判定処理は終了する。

ステップS 4 0 8では、音声成分強調部2 1 0において、音を強調する方向を更新する。

ステップS 4 1 0では、雑音成分強調部2 1 2において、フィルタ行列の最適化処理を行って、フィルタ行列の更新を行う。

なお本実施例の音声認識装置1 0 0では、ステップS 4 0 8とステップS 4 1 0は、ステップS 4 0 6で更新が必要と判定された後に、同時に並行して実施される。ステップS 4 0 8とステップS 4 1 0の両方の処理が終了した後、図4の更新判定処理は終了する。

【0 0 6 0】

本実施例の音声認識装置1 0 0によれば、雑音成分を推定するにあたり、対話者Vの方向からの音を抑圧するのではなく、ICAに基づいた信号の分離によって、雑音成分を推定する。このような手法を用いることによって、マイクロホンの感度特性のばらつきの影響を受けることなく、雑音成分を推定することができる。これによって、雑音成分を正確に除去した音声メル周波数スペクトルを得ることができる。言語認識部2 2 2における誤認識を防ぐことができる。

【0 0 6 1】

本実施例の音声認識装置1 0 0によれば、雑音成分推定部2 1 2におけるフィルタ行列を、常に更新し続けるのではなく、対話者Vの方向が所定の角度以上変化した場合に更新する構成としている。フィルタ行列の更新を頻繁に行う場合、処理の負荷は非常に重いものになってしまうが、本実施例のように必要な時にのみフィルタ行列を更新する構成とすることによって、処理の負荷を大幅に軽減することができる。音声認識処理に要する時間を短時間にするすることができる。

【0 0 6 2】

上記の実施例では、対話者Vの方向が所定の角度以上変化した場合に、更新判定部2 0 4が更新を指示する例を説明したが、更新の判定基準はこれに限らない。例えば、所定の時間（例えば1 0 s）が経過するごとに、更新を指示する構成としてもよい。あるいは、音声認識装置1 0 0が対話者を新たに検出するごとに、更新を指示する構成としてもよい。

【0 0 6 3】

また上記とは異なり、雑音成分推定部2 1 2におけるフィルタ行列の更新について、音声認識装置1 0 0および対話者Vの周囲の環境が変化した場合に、更新判定部2 0 4で更新を指示する構成としてもよい。この場合、画像認識部2 0 2において、右カメラ1 0 4と左カメラ1 0 6から出力される画像データから、対話者Vの輪郭を抽出して、対話者Vの輪郭から対話者Vの映像と周囲の環境の映像を識別し、周囲の環境の映像の経時的変化

を特定することで、周囲の環境の変化を認識することができる。この場合、更新判定部 204 は、画像認識部 202 から周囲の環境が変化したことを通知されると、雑音成分推定部 212 へフィルタ行列の更新を指示する。

【0064】

上記の実施例では、音声成分強調部 210 において音を強調する方向の更新と、雑音成分推定部 212 におけるフィルタ行列の更新を同時に行う例を説明したが、これらは別々のタイミングで更新する構成としてもよい。特に、音声成分強調部 210 において音を強調する方向の更新は、雑音成分推定部 212 におけるフィルタ行列の更新に比べて処理の負荷が軽いので、音声成分強調部 210 において音を強調する方向は、対話者 V の方向が変化する度に更新する構成としてもよい。

10

【0065】

上記の実施例では、音声成分強調部 210 で算出された基本周波数スペクトルと、雑音成分推定部 212 で算出された雑音周波数スペクトルのそれぞれについて、先ずメル周波数変換部 214 および 216 で別個にメル周波数スペクトルへの変換を行い、その後スペクトル減算部 218 でスペクトル・サブトラクションを行っている。上記とは異なり、先ず音声成分強調部 210 で算出された基本周波数スペクトルと、雑音成分推定部 212 で算出された雑音周波数スペクトルについてスペクトル・サブトラクションを行い、その後メル周波数スペクトルへの変換を行う構成としてもよい。

【0066】

以上、本発明の具体例を詳細に説明したが、これらは例示にすぎず、特許請求の範囲を限定するものではない。特許請求の範囲に記載の技術には、以上に例示した具体例を様々な変形、変更したものが含まれる。

20

また、本明細書または図面に説明した技術要素は、単独であるいは各種の組み合わせによって技術的有用性を発揮するものであり、出願時請求項記載の組み合わせに限定されるものではない。また、本明細書または図面に例示した技術は複数目的を同時に達成するものであり、そのうちの一つの目的を達成すること自体で技術的有用性を持つものである。

【図面の簡単な説明】

【0067】

【図 1】図 1 は音声認識装置 100 の外観を示す図である。

【図 2】図 2 はコントローラ 114 の構成を模式的に示す図である。

30

【図 3】図 3 は音声認識装置 100 が実施する音声認識処理のフローチャートである。

【図 4】図 4 は音声認識装置 100 が実施する更新判定処理のフローチャートである。

【図 5】図 5 は独立成分分析 (ICA) の概要を説明する図である。

【図 6】図 6 はスペクトル・サブトラクション (SS) 法の概要を説明する図である。

【図 7】図 7 は SS 法における指向特性を示す図である。

【図 8】図 8 は SS 法におけるマイクロホンの感度特性と指向特性の関係を示す図である。

。

【図 9】図 9 は SS 法におけるマイクロホンの感度特性と指向特性の関係を示す図である。

。

【図 10】図 10 は周波数変換部 208 の処理の概要を説明する図である。

40

【符号の説明】

【0068】

100 : 音声認識装置

102 : 頭部

104 : 右カメラ

106 : 左カメラ

108 : 胴体部

110 : 集音部

112 a、112 b、112 c、112 d、112 e、112 f : マイクロホン

114 : コントローラ

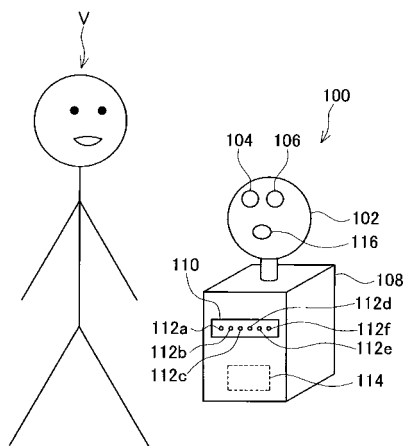
50

116 : スピーカ
 202 : 画像認識部
 204 : 更新判定部
 206 : A / D 変換部
 208 : 周波数変換部
 210 : 音声成分強調部
 212 : 雑音成分推定部
 214、216 : メル周波数変換部
 218 : スペクトル減算部
 220 : 特徴量計算部
 222 : 言語認識部
 224 : 応答制御部
 226 : D / A 変換部
 502、504 : 音源
 506、508 : マイクロホン
 510 : フィルタ行列
 512 : フィルタ最適化手段
 602 : 音源
 604、606 : マイクロホン
 610、612 : 経路
 616、618 : 波面
 620、622 : 点
 624、626 : 音信号
 702、704、810、812、910、912 : 指向特性の分布
 1002 : 音データ

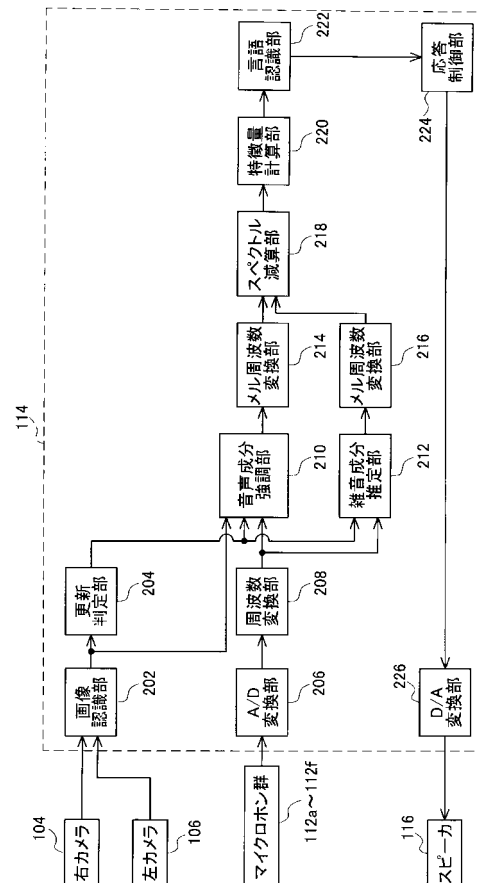
10

20

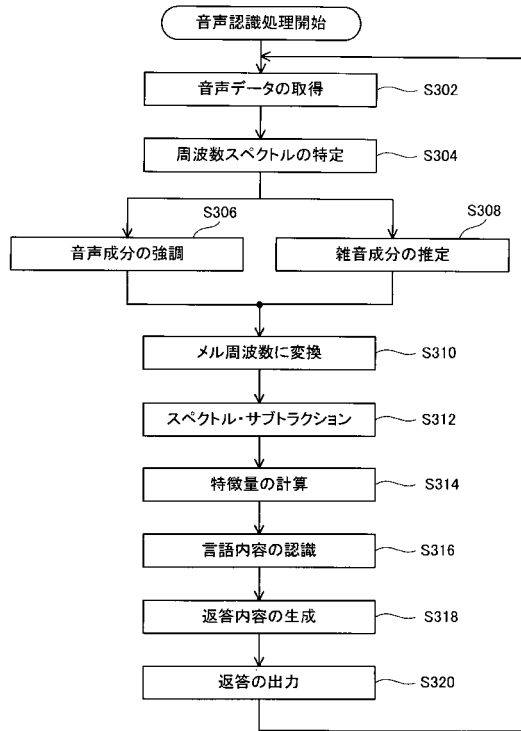
【図 1】



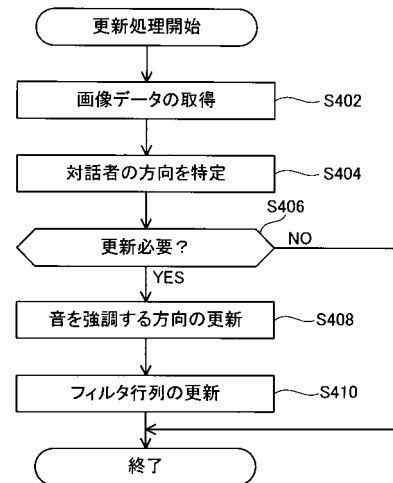
【図 2】



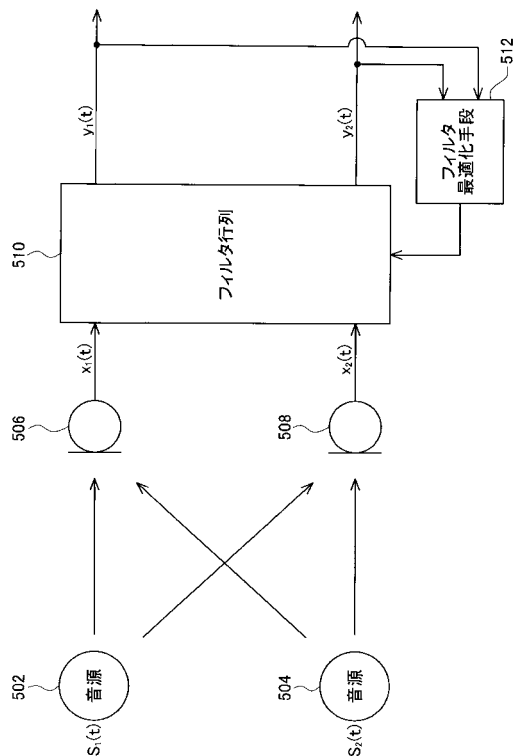
【図 3】



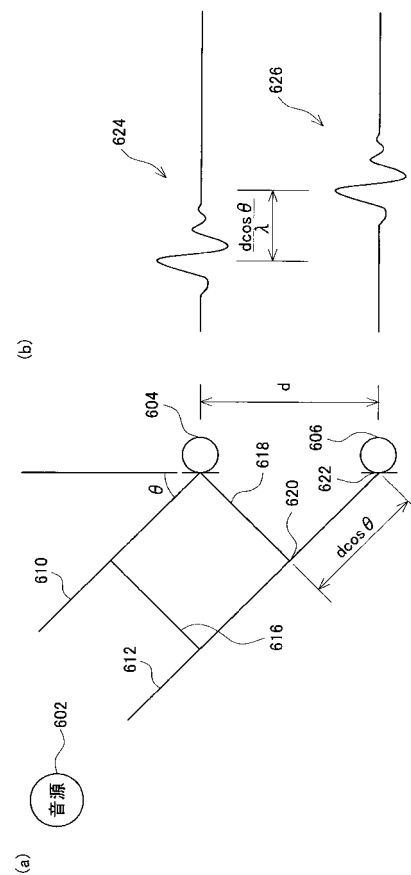
【図 4】



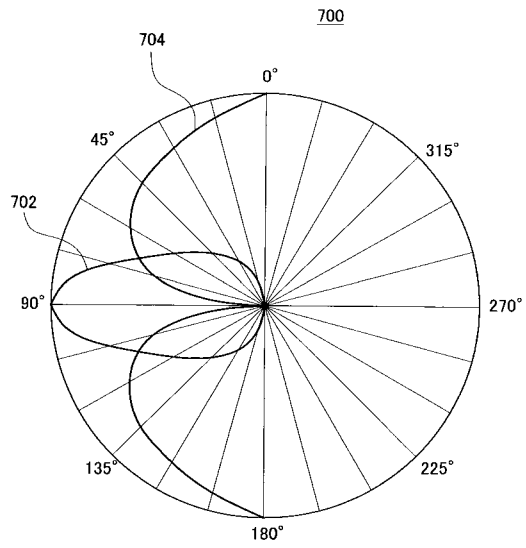
【図 5】



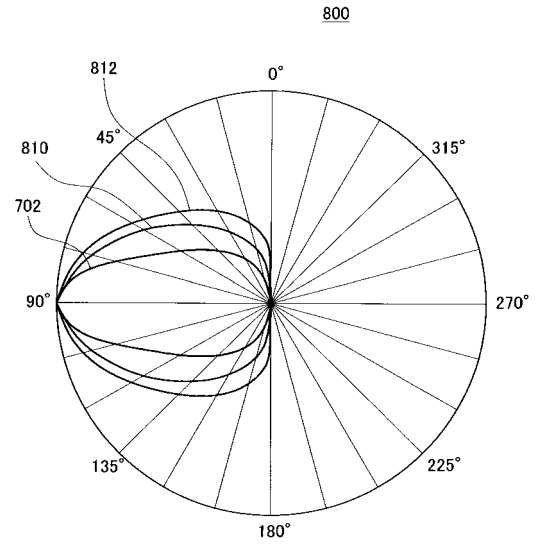
【図 6】



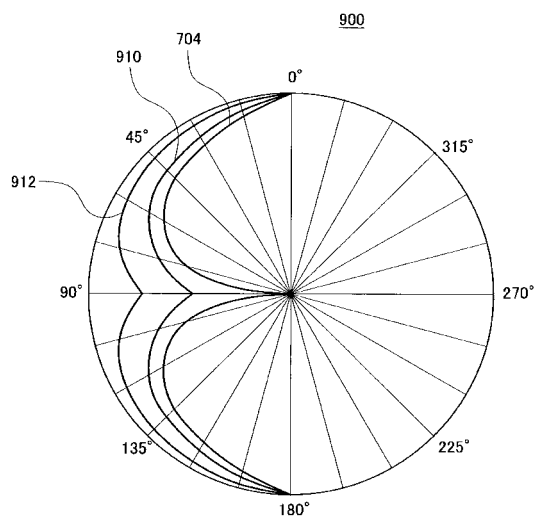
【図 7】



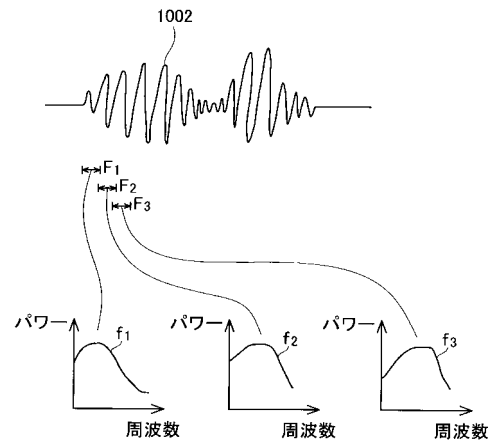
【図 8】



【図 9】



【図 10】



フロントページの続き

(51)Int.Cl. F I
G 0 6 T 7/60 1 5 0 P

(72)発明者 高橋 祐
奈良県生駒市高山町 8 9 1 6 - 5 国立大学法人 奈良先端科学技術大学院大学内

(72)発明者 渡部 生聖
愛知県豊田市トヨタ町 1 番地 トヨタ自動車株式会社内

審査官 毛利 太郎

(56)参考文献 特開 2 0 0 4 - 0 6 9 7 7 2 (J P , A)
特開 2 0 0 3 - 0 6 6 9 8 6 (J P , A)
特開 2 0 0 4 - 1 3 3 4 0 3 (J P , A)
特開 2 0 0 5 - 3 0 8 7 7 1 (J P , A)
特開 2 0 0 5 - 0 3 1 2 5 8 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)
G 1 0 L 1 1 / 0 0 - 2 1 / 0 6
G 0 6 T 7 / 6 0