

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4516527号  
(P4516527)

(45) 発行日 平成22年8月4日(2010.8.4)

(24) 登録日 平成22年5月21日(2010.5.21)

(51) Int.Cl. F I  
**G 1 0 L 15/06 (2006.01)** G 1 0 L 15/06 3 1 0 T  
**G 1 0 L 15/28 (2006.01)** G 1 0 L 15/06 4 0 0 V  
 G 1 0 L 15/28 4 0 0

請求項の数 7 (全 29 頁)

(21) 出願番号	特願2005-515466 (P2005-515466)	(73) 特許権者	000005326 本田技研工業株式会社 東京都港区南青山二丁目1番1号
(86) (22) 出願日	平成16年11月12日(2004.11.12)	(74) 代理人	100064414 弁理士 磯野 道造
(86) 国際出願番号	PCT/JP2004/016883	(72) 発明者	中臺 一博 埼玉県和光市中央1丁目4-1
(87) 国際公開番号	W02005/048239	(72) 発明者	辻野 広司 埼玉県和光市中央1丁目4-1
(87) 国際公開日	平成17年5月26日(2005.5.26)	(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 フォルム東洞院三条 1102号
審査請求日	平成19年7月11日(2007.7.11)	(72) 発明者	山本 俊一 埼玉県和光市本町8-1
(31) 優先権主張番号	特願2003-383072 (P2003-383072)		
(32) 優先日	平成15年11月12日(2003.11.12)		
(33) 優先権主張国	日本国(JP)		

最終頁に続く

(54) 【発明の名称】 音声認識装置

(57) 【特許請求の範囲】

【請求項1】

複数のマイクが検出した音響信号から、音声認識して文字情報に変換する音声認識装置であって、

前記複数のマイクが検出した音響信号に基づき、特定の話者の音源方向を特定する音源定位部と、

前記複数のマイクが検出した1つ以上の音響信号に基づき、その音響信号に含まれる音声信号の特徴を抽出する特徴抽出部と、

断続的な複数の方向に対応した方向依存音響モデルを記憶した音響モデル記憶部と、

前記音源定位部が特定した音源方向の音響モデルを、前記音響モデル記憶部の方向依存音響モデルと当該方向依存音響モデル毎に設定された重みとを内積して合成して、前記音響モデル記憶部へ記憶させるパラメータ合成部を備える音響モデル合成部と、

前記音響モデル合成部が合成した音響モデルを使用して、前記特徴抽出部が抽出した特徴について音声認識を行い、文字情報に変換する音声認識部と、を備え、

前記パラメータ合成部は、

前記音源が正面にあるときの重みを定める関数を学習により設定し、前記音源方向に対応する音響モデルを合成する際、前記音源が正面にあるときの重みを定める関数を前記音源方向に移動した関数を求め、当該移動した関数を参照して重みを設定し、

前記学習として、前記音源が正面にあるときの重み初期値が予め設定され、当該重み初期値を用いて合成した音響モデルで前記音素列を認識させ、正解を出した前記方向依存音

響モデルの重みを増加させ、正解を出さなかった前記方向依存音響モデルの重みを減少させて更新する試行を行うと共に更新した前記方向依存音響モデルの重みを用いて前記試行を所定の回数繰り返すことで、前記更新した方向依存音響モデルの重みを、前記音源が正面にあるときの重みを定める関数として設定することを特徴とする音声認識装置。

【請求項 2】

複数のマイクが検出した音響信号から、特定の話者の音声を認識して文字情報に変換する音声認識装置であって、

前記複数のマイクが検出した音響信号に基づき、前記特定の話者の音源方向を特定する音源定位部と、

前記音源定位部が特定した音源方向に基づき、前記特定の話者の音声信号を前記音響信号から分離する音源分離部と、

前記音源分離部が分離した音声信号の特徴を抽出する特徴抽出部と、

断続的な複数の方向に対応した方向依存音響モデルを記憶した音響モデル記憶部と、

前記音源定位部が特定した音源方向の音響モデルを、前記音響モデル記憶部の方向依存音響モデルと当該方向依存音響モデル毎に設定された重みとを内積して合成して、前記音響モデル記憶部へ記憶させるパラメータ合成部を備える音響モデル合成部と、

前記音響モデル合成部が合成した音響モデルを使用して、前記特徴抽出部が抽出した特徴について音声認識を行い、文字情報に変換する音声認識部と、を備え、

前記パラメータ合成部は、

前記音源が正面にあるときの重みを定める関数を学習により設定し、前記音源方向に対応する音響モデルを合成する際、前記音源が正面にあるときの重みを定める関数を前記音源方向に移動した関数を求め、当該移動した関数を参照して重みを設定し、

前記学習として、前記音源が正面にあるときの重み初期値が予め設定され、当該重み初期値を用いて合成した音響モデルで前記音素列を認識させ、正解を出した前記方向依存音響モデルの重みを増加させ、正解を出さなかった前記方向依存音響モデルの重みを減少させて更新する試行を行うと共に更新した前記方向依存音響モデルの重みを用いて前記試行を所定の回数繰り返すことで、前記更新した方向依存音響モデルの重みを、前記音源が正面にあるときの重みを定める関数として設定することを特徴とする音声認識装置。

【請求項 3】

前記音源分離部は、前記音源定位部が特定した音源方向が、前記複数のマイクの配置により決定される正面に近い場合には、狭い方向帯域の音声を分離し、正面から離れると広い方向帯域の音声を分離するアクティブ方向通過型フィルタを用いて音声分離を行うよう構成されたことを特徴とする請求項 2 に記載の音声認識装置。

【請求項 4】

前記音源定位部は、前記マイクが検出した音響信号を周波数分析した後、調波構造を抽出し、複数のマイクから抽出された調波構造の音圧差と位相差とを求め、この音圧差と位相差のそれぞれから音源方向の確からしさを求め、最も確からしい方向を音源方向と判断するよう構成されたことを特徴とする請求項 1 または請求項 2 に記載の音声認識装置。

【請求項 5】

前記音源定位部は、前記複数のマイクから検出された音響信号の音圧差と位相差を用いて前記特定の話者の音源方向を特定するために、前記マイクが設けられる部材の表面で散乱する音響信号を音源方向ごとにモデル化した散乱理論を用いることを特徴とする請求項 1、請求項 2 または請求項 4 のいずれか 1 項に記載の音声認識装置。

【請求項 6】

前記音響モデル合成部は、前記音響モデル記憶部の方向依存音響モデルの重み付き線形和により前記音源方向の音響モデルを合成するよう構成され、

前記線形和に使用する重みが、学習により決定されたことを特徴とする請求項 1 から請求項 5 のいずれか 1 項に記載の音声認識装置。

【請求項 7】

前記話者を特定する話者同定部をさらに備え、

10

20

30

40

50

前記音響モデル記憶部は、前記話者ごとに方向依存音響モデルを有し、

前記音響モデル合成部は、前記話者同定部が特定した話者の方向依存音響モデルと、前記音源定位部が特定した音源方向とに基づき、前記音源方向の音響モデルを前記音響モデル記憶部の方向依存音響モデルに基づいて求め、前記音響モデル記憶部へ記憶させるよう構成されたことを特徴とする請求項 1 から請求項 6 のいずれか 1 項に記載の音声認識装置

。【発明の詳細な説明】

【技術分野】

【0001】

本発明は、音声認識装置に関し、詳しくは、話者や、音声認識装置を備えた移動体が移動しても高い精度で音声を認識可能な音声認識装置に関する。 10

【背景技術】

【0002】

近年、音声認識技術は、実用化の域に入ってきており、情報の音声入力などに利用され始めている。一方、ロボットの研究開発も盛んとなっており、音声認識技術は、ロボットを実用化するための一つのキー技術ともなっている。すなわち、ロボットと人間との知的なソーシャルインタラクションを行うためには、人間の言葉をロボットが理解する必要があるため、音声認識の精度が重要となっている。

【0003】

ところが、実際に人とのコミュニケーションを行うためには、実験室において口元に設置したマイクで音声を入力して行う音声認識とは異なるいくつかの問題がある。 20

例えば、実際の環境には様々な雑音があり、雑音の中から必要な音声信号を抽出しなければ音声認識をすることができない。また、話者が複数存在する場合にも、同様に認識の対象とする話者の音声のみを抽出する必要がある。また、音声認識においては、一般に隠れマルコフモデル(HMM: Hidden Markov Model)というモデルを利用して内容を特定するが、話者の位置(音響認識装置のマイクを基準とした方向)が異なると、話者の声の聞こえ方も異なることから、認識率に影響を及ぼすという問題がある。

【0004】

このようなことから、本発明者を含む研究グループでは、アクティブオーディションにより複数の音源の定位・分離・認識を行う技術を発表している(非特許文献1参照)。 30

この技術は、人間の耳に相当する位置に2つのマイクを配置し、複数の話者が同時に発話した場合に、一人の発した単語を認識する技術である。詳しくは、2つのマイクから入力された音響信号から、話者の位置を定位し、各話者の音声を分離した上で、音声認識する。この認識の際、移動体(音声認識装置を備えたロボット等)から見て-90°から90°まで10°おきの方向に対する各話者の音響モデルを予め作成しておく。そして、音声の認識時には、それらの音響モデルを用いて並列に認識プロセスを実行する。

【非特許文献1】 A Humanoid Listens to three simultaneous talkers by Integrating Active Audition and Face Recognition Kazuhiro Nakadai, et al., IJCAI-03 Workshop on Issues in Designing Physical Agents for Dynamic Real-Time Environments: World Modeling, Planning, Learning and Communicating, PP 117-124 40

【発明の開示】

【0005】

しかしながら、前記した従来技術では、話者や移動体が移動する場合には、その都度移動体に対する話者の位置が変化するため、予め用意された音響モデルの方向と異なる方向に話者が位置すると、認識率が低下するという問題があった。 50

本発明は、このような背景に鑑みてなされたもので、話者や、移動体が移動しても高い精度で認識可能な音声認識装置を提供することを課題とする。

【0006】

前記課題を解決するため、本発明の音声認識装置は、複数のマイクが検出した音響信号から、音声を認識して文字情報に変換する音声認識装置であって、前記複数のマイクが検出した音響信号に基づき、特定の話者の音源方向を特定する音源定位部と、前記複数のマイクが検出した1つ以上の音響信号に基づき、その音響信号に含まれる音声信号の特徴を抽出する特徴抽出部と、断続的な複数の方向に対応した方向依存音響モデルを記憶した音響モデル記憶部と、前記音源定位部が特定した音源方向の音響モデルを、前記音響モデル記憶部の方向依存音響モデルと当該方向依存音響モデル毎に設定された重みとを内積して合成して、前記音響モデル記憶部へ記憶させるパラメータ合成部を備える音響モデル合成部と、前記音響モデル合成部が合成した音響モデルを使用して、前記特徴抽出部が抽出した特徴について音声認識を行い、文字情報に変換する音声認識部と、を備え、前記パラメータ合成部は、前記音源が正面にあるときの重みを定める関数を学習により設定し、前記音源方向に対応する音響モデルを合成する際、前記音源が正面にあるときの重みを定める関数を前記音源方向に移動した関数を求め、当該移動した関数を参照して重みを設定し、前記学習として、前記音源が正面にあるときの重み初期値が予め設定され、当該重み初期値を用いて合成した音響モデルで前記音素列を認識させ、正解を出した前記方向依存音響モデルの重みを増加させ、正解が出さなかった前記方向依存音響モデルの重みを減少させて更新する試行を行うと共に更新した前記方向依存音響モデルの重みを用いて前記試行を所定の回数繰り返すことで、前記更新した方向依存音響モデルの重みを、前記音源が正面にあるときの重みを定める関数として設定することを特徴とする。

10

20

【0007】

このような音声認識装置によれば、音源定位部が音源方向を特定し、音響モデル合成部は、音源方向と、方向依存音響モデルとに基づき、その方向に適した音響モデルを合成し、音声認識部がこの音響モデルを使用して音声認識を行う。

【0008】

また、前記した音声認識装置においては、音源定位部が特定した音源方向に基づき、前記特定の話者の音声信号を前記音響信号から分離する音源分離部を備え、音源分離部が分離した音声信号に基づき、特徴抽出部が音声信号の特徴を抽出するように構成してもよい。

30

【0009】

このような音声認識装置によれば、音源定位部が音源方向を特定し、音源分離部は、音源定位部が特定した音源方向の音声のみを分離する。そして、音響モデル合成部は、音源方向と、方向依存音響モデルとに基づき、その方向に適した音響モデルを合成し、音声認識部がこの音響モデルを使用して音声認識を行う。

なお、音源分離部が出力する音声信号というのは、音声としての意味を持つ情報であればよく、音声のアナログ信号そのものに限らず、デジタル化、符号化した信号や、周波数分析したスペクトルのデータを含む。

【0010】

また、前記した音声認識装置では、前記音源定位部は、前記マイクが検出した音響信号を周波数分析した後、調波構造を抽出し、複数のマイクから抽出された調波構造の音圧差と位相差とを求め、この音圧差と位相差のそれぞれから音源方向の確からしさを求め、最も確からしい方向を音源方向と判断するよう構成することができる。

40

【0011】

また、前記音源定位部は、前記複数のマイクから検出された音響信号の音圧差と位相差を用いて前記特定の話者の音源方向を特定するために、ロボットの頭部などの前記マイクが設けられる部材の表面で散乱する音響信号を音源方向ごとにモデル化した散乱理論を用いることができる。

【0012】

50

さらに、前記した音声認識装置では、前記音源分離部は、前記音源定位部が特定した音源方向が、前記複数のマイクの配置により決定される正面に近い場合には、狭い方向帯域の音声を分離し、正面から離れると広い方向帯域の音声を分離するアクティブ方向通過型フィルタを用いて音声分離を行うよう構成されるのが好ましい。

【0013】

また、前記した音声認識装置では、前記音響モデル合成部は、前記音響モデル記憶部の方向依存音響モデルの重み付き線形和により前記音源方向の音響モデルを合成するよう構成され、前記線形和に使用する重みが、学習により決定されるのが好ましい。

【0014】

また、前記した音声認識装置では、前記話者を特定する話者同定部をさらに備え、前記音響モデル記憶部は、前記話者ごとに方向依存音響モデルを有し、前記音響モデル合成部は、前記話者同定部が特定した話者の方向依存音響モデルと、前記音源定位部が特定した音源方向に基づき、前記音源方向の音響モデルを前記音響モデル記憶部の方向依存音響モデルに基づいて求め、前記音響モデル記憶部へ記憶させるよう構成されるのが好ましい。

10

【0015】

また、前記特徴抽出部で抽出された特徴、または前記音源分離部が分離した音声信号について、予め用意した雛形と比較し、前記雛形との違いが予め設定した閾値より大きい領域、例えば周波数領域や、サブバンドを同定し、同定された領域については、その特徴としての信頼性が低いことを示す指標を前記音声認識部へ出力するマスキング部をさらに備えるのが望ましい。

20

【図面の簡単な説明】

【0018】

【図1】本発明の実施形態に係る音声認識装置のブロック図である。

【図2】音源定位部の一例を示すブロック図である。

【図3】音源定位部の動作を説明する図である。

【図4】音源定位部の動作を説明する図である。

【図5】聴覚エピソード幾何を説明する図である。

【図6】位相差と周波数  $f$  の関係を示すグラフである。

【図7】頭部伝達関数の一例を示すグラフである。

30

【図8】音源分離部の一例を示すブロック図である。

【図9】通過帯域関数の一例を示すグラフである。

【図10】サブバンド選択部の動作を説明する図である。

【図11】通過帯域の一例を図示した平面図である。

【図12】(a) および (b) は、ともに特徴抽出部の一例を示すブロック図である。

【図13】音響モデル合成部の一例を示すブロック図である。

【図14】方向依存音響モデルの認識単位とサブモデルを示した図である。

【図15】パラメータ合成部の動作を説明する図である。

【図16】(a) および (b) は、ともに重み  $W_n$  の一例を示すグラフである。

【図17】重み  $W$  の学習方法を説明する図である。

40

【図18】第2実施形態に係る音声認識装置のブロック図である。

【図19】音響の入力距離差を示す図である。

【図20】第3実施形態に係る音声認識装置のブロック図である。

【図21】ストリーム追跡部のブロック図である。

【図22】音源方向の履歴を図示したグラフである。

【発明を実施するための最良の形態】

【0019】

[第1実施形態]

次に、本発明の実施形態について、適宜図面を参照しながら詳細に説明する。図1は、本発明の実施形態に係る音声認識装置のブロック図である。

50

図1に示すように、実施形態に係る音声認識装置1は、2つのマイク $M_R$ 、 $M_L$ と、マイク $M_R$ 、 $M_L$ が検出した音響信号から、話者(音源)の位置を特定する音源定位部10と、音源定位部10が特定した音源方向及び音源定位部10で求めたスペクトルに基づいて、特定の方向の音源から来る音響を分離する音源分離部20と、複数の方向についての音響モデルを記憶した音響モデル記憶部49と、音響モデル記憶部49内の音響モデル及び音源定位部10が特定した音源方向に基づいて、その音源方向の音響モデルを合成する音響モデル合成部40と、音源分離部20が分離した特定音源のスペクトルから音響の特徴を抽出する特徴抽出部30と、音響モデル合成部40が合成した音響モデルと、特徴抽出部30が抽出した音響の特徴に基づき音声認識を行う音声認識部50とを備える。これらのうち、音源分離部20は、任意的に用いられる。

10

本発明では、音響モデル合成部40が生成した、音源の方向に適した音響モデルを利用して音声認識部50が音声認識を行うため、高い認識率が実現される。

#### 【0020】

次に、実施形態に係る音声認識装置1の構成要素であるマイク $M_R$ 、 $M_L$ 、音源定位部10、音源分離部20、特徴抽出部30、音響モデル合成部40、及び音声認識部50についてそれぞれ説明する。

#### 【0021】

##### 《マイク $M_R$ 、 $M_L$ 》

マイク $M_R$ 、 $M_L$ は、音を検出して電気信号(音響信号)として出力する一般的なマイクである。本実施形態では、2つとしているが、複数であれば幾つでもよく、例えば3つ、4つを使用しても構わない。マイク $M_R$ 、 $M_L$ は、例えば、移動体であるロボットRBの両耳の部分に設けられる。

20

マイク $M_R$ 、 $M_L$ の配置は、音響信号を集音するための一般的な音声認識装置1の正面を決定する。すなわち、マイク $M_R$ 、 $M_L$ の集音方向のベクトルの和の方向が音声認識装置1の正面となる。図1に示すように、ロボットRBの頭の左右両脇にマイク $M_R$ 、 $M_L$ が1つずつ設けられていれば、ロボットRBの正面が音声認識装置1の正面となる。

#### 【0022】

##### 《音源定位部10》

図2は、音源定位部の一例を示すブロック図であり、図3及び図4は、音源定位部の動作を説明する図である。

30

音源定位部10は、2つのマイク $M_R$ 、 $M_L$ から入力された2つの音響信号から、各話者 $HM_j$ (図3では、 $HM_1$ 、 $HM_2$ )の音源方向を定位する。音源定位方法は、マイク $M_R$ 、 $M_L$ に入力された音響信号の位相差を利用する方法、ロボットRBの頭部伝達関数を用いて推定する方法、右と左のマイク $M_R$ 、 $M_L$ から入力された信号の相互相関をとる方法などがあり、それぞれ精度を上げるため、種々の改良が加えられているが、ここでは、本発明者が改良した手法を例にして説明する。

#### 【0023】

音源定位部10は、図2に示すように、周波数分析部11、ピーク抽出部12、調波構造抽出部13、IPD計算部14、IID計算部15、聴覚エピソード幾何仮説データ16、確信度計算部17、及び確信度統合部18を備える。

40

これらの各部を、図3及び図4を参照しながら説明する。場面として、ロボットRBに対し、2人の話者 $HM_1$ 、 $HM_2$ が同時に話しかける場合で説明する。

#### 【0024】

##### 周波数分析部11

周波数分析部11は、ロボットRBが備える左右のマイク $M_R$ 、 $M_L$ が検出した左右の音響信号 $CR_1$ 、 $CL_1$ から、微小時間 $t$ の時間長の信号区間を切り出し、左右のチャンネルごとにFFT(高速フーリエ変換)により周波数分析を行う。

例えば、右のマイク $M_R$ からの音響信号 $CR_1$ より得られる分析結果がスペクトル $CR_2$ であり、左のマイク $M_L$ からの音響信号 $CL_1$ より得られる分析結果がスペクトル $CL_2$ である。

50

なお、周波数分析は、バンドパスフィルタなど、他の手法を用いることもできる。

【 0 0 2 5 】

ピーク抽出部 1 2

ピーク抽出部 1 2 は、スペクトル C R 2 , C L 2 から左右のチャンネルごとに一連のピークを抽出する。ピークの抽出は、スペクトルのローカルピークをそのまま抽出するか、スペクトラルサブトラクション法に基づいた方法 ( S . F . Boll , A s p e c t r a l s u b t r a c t i o n a l g o r i t h m f o r s u p p r e s s i o n o f a c o u s t i c n o i s e i n s p e e c h , P r o c e e d i n g s o f 1 9 7 9 I n t e r n a t i o n a l c o n f e r e n c e o n A c o u s t i c s , S p e e c h , a n d s i g n a l P r o c e s s i n g ( I C A S S P - 7 9 ) 参照 ) で行う。後者の方法は、スペクトルからピークを抽出し、これをスペクトルから減算し、残差スペクトルを生成する。そして、その残差スペクトルからピークが見つからなくなるまでピーク抽出の処理を繰り返す。

10

前記スペクトル C R 2 , C L 2 に対し、ピークの抽出を行うと、例えばピークスペクトル C R 3 , C L 3 のようにピークを構成するサブバンドの信号のみが抽出される。

【 0 0 2 6 】

調波構造抽出部 1 3

調波構造抽出部 1 3 は、音源が有する調波構造に基づき、左右のチャンネルごとに特定の調波構造を有するピークをグループにする。例えば、人の声であれば、特定の人声は、基本周波数の音と、基本周波数の倍音とからなるが、人により基本周波数が微妙に異なるので、その周波数の差により、複数の人声をグループ分けすることができる。調波構造に基づいて同じグループに分けられたピークは、同じ音源から発せられた信号と推定できる。例えば、複数 ( J 人 ) の話者が同時に話していれば、複数 ( J 個 ) の調波構造が抽出される。

20

【 0 0 2 7 】

図 3 においては、ピークスペクトル C R 3 , C L 3 の、ピーク P 1 , P 3 , P 5 を一つのグループにして調波構造 C R 4 1 , C L 4 1 とし、ピーク P 2 , P 4 , P 6 を一つのグループにして調波構造 C R 4 2 , C L 4 2 としている。

【 0 0 2 8 】

I P D 計算部 1 4

I P D 計算部 1 4 は、調波構造抽出部 1 3 が抽出した調波構造 C R 4 1 , C R 4 2 , C L 4 1 , C L 4 2 のスペクトルから、I P D ( 両耳間位相差 ) を計算する部分である。

30

I P D 計算部 1 4 は、話者 H M j に対応する調波構造 ( 例えば、調波構造 C R 4 1 ) に含まれているピーク周波数の集合を  $\{ f_k \mid k = 0 \dots K - 1 \}$  としたとき、各  $f_k$  に対応するスペクトルのサブバンドを、右と左の両チャンネル ( 例えば、調波構造 C R 4 1 と調波構造 C L 4 1 ) から選択し、次式 ( 1 ) により I P D (  $f_k$  ) を計算する。調波構造 C R 4 1 と調波構造 C L 4 1 から計算した I P D (  $f_k$  ) は、例えば、図 4 に示す両耳間位相差 C 5 1 のようになる。ここで、(  $f_k$  ) は、ある調波構造に含まれるある倍音  $f_k$  の I P D であり、K は、その調波構造に含まれる倍音の数を示す。

【 0 0 2 9 】

40

【 数 1 】

$$\Delta \phi(f_k) = \arctan\left(\frac{\Im[S_r(f_k)]}{\Re[S_r(f_k)]}\right) - \arctan\left(\frac{\Im[S_l(f_k)]}{\Re[S_l(f_k)]}\right) \dots (1)$$

但し、

- (  $f_k$  ) :  $f_k$  の I P D ( 両耳間位相差 )
- J [ S<sub>r</sub> (  $f_k$  ) ] : 右の入力信号のピーク  $f_k$  のスペクトル虚部
- R [ S<sub>r</sub> (  $f_k$  ) ] : 右の入力信号のピーク  $f_k$  のスペクトル実部
- J [ S<sub>l</sub> (  $f_k$  ) ] : 左の入力信号のピーク  $f_k$  のスペクトル虚部

50

$R[S_L(f_k)]$  : 左の入力信号のピーク  $f_k$  のスペクトル実部

【0030】

IID計算部15

IID計算部15は、各調波構造にある各倍音について、左のマイク  $M_L$  から入力された音の音圧と、右のマイク  $M_R$  から入力された音の音圧との差（両耳間音圧差）を計算する部分である。

IID計算部15は、話者  $HM_j$  に対応する調波構造（例えば、調波構造  $CR41$  ,  $CL41$ ）に含まれているピーク周波数  $f_k$  の倍音に対応するスペクトルのサブバンドを、右と左の両チャンネル（例えば、調波構造  $CR41$  と調波構造  $CL41$ ）から選択し、次式（2）により  $IID(f_k)$  を計算する。調波構造  $CR41$  と調波構造  $CL41$  から計算した  $IID(f_k)$  は、例えば図4に示す両耳間音圧差  $C61$  のようになる。

10

【0031】

【数2】

$$\Delta p(f_k) = p_r(f_k) - p_l(f_k) \quad \dots (2)$$

但し、

$(f_k)$  :  $f_k$  の IID（両耳間音圧差）

$p_r(f_k)$  : 右の入力信号のピーク  $f_k$  のパワー

$p_l(f_k)$  : 左の入力信号のピーク  $f_k$  のパワー

$$p_r(f_k) = 10 \log_{10} (J[S_r(f_k)]^2 + R[S_r(f_k)]^2) \quad 20$$

$$p_l(f_k) = 10 \log_{10} (J[S_l(f_k)]^2 + R[S_l(f_k)]^2)$$

【0032】

聴覚エピソード幾何仮説データ16

聴覚エピソード幾何仮説データ16は、図5に示すように、ロボット  $RB$  の頭部を想定した球体を上から見たときに、音源  $S$  と、ロボット  $RB$  の両耳のマイク  $M_R$  ,  $M_L$  との距離差から生じる時間差に基づき想定される位相差のデータである。

聴覚エピソード幾何により、位相差は、次式（3）により求められる。ここでは、頭部形状を球と仮定している。

【0033】

【数3】

30

$$\Delta \phi = \frac{2\pi f}{v} \times r(\theta + \sin \theta) \quad \dots (3)$$

【0034】

ここで、 $\Delta \phi$  は両耳間位相差（IPD）、 $v$  は音速、 $f$  は周波数、 $r$  は両耳間の距離  $2r$  から求まる値、 $\theta$  は音源方向を示す。

式（3）により、各音源方向より発せられた音響信号の周波数  $f$  と位相差  $\Delta \phi$  の関係は、図6のようになる。

【0035】

確信度計算部17

40

確信度計算部17は、IPD及びIIDのそれぞれの確信度を計算する。

- IPD確信度 -

IPDの確信度は、話者  $HM_j$  に対応する調波構造（例えば、調波構造  $CR41$  ,  $CL41$ ）が含んでいる倍音  $f_k$  がどの方向から来ているらしいかを  $\Delta \phi$  の関数として求め、これを確率関数にあてはめる。

まず、 $f_k$  の IPD の仮説（予想値）を次式（4）に基づき計算する。

【0036】



【数4】

$$\Delta \phi_h(\theta, f_k) = \frac{2\pi f_k}{v} \times r(\theta + \sin \theta) \quad \dots (4)$$

【0037】

$\phi_h(\theta, f_k)$  は、ある調波構造内の  $k$  番目の倍音  $f_k$  に対して音源方向が  $\theta$  の場合の IPD の仮説 (予想値) を示す。IPD の仮説は、例えば音源方向  $\theta$  を、 $\pm 90^\circ$  の範囲で  $5^\circ$  おきに変化させて計 37 個の仮説を計算する。もっとも、より細かい角度ごとに計算しても、より大まかな角度ごとに計算してもかまわない。

次に、次式 (5) により、 $\phi_h(\theta, f_k)$  と  $\phi(f_k)$  の差を求め、すべてのピーク  $f_k$  について合計する。この差は、仮説と入力との距離を表し、 $\theta$  が話者のいる方向に近いと小さく、遠いと大きくなる。

【0038】

【数5】

$$d(\theta) = \frac{1}{K} \sum_{k=0}^{K-1} \frac{(\Delta \phi_h(\theta, f_k) - \Delta \phi(f_k))^2}{f_k} \quad \dots (5)$$

【0039】

得られた  $d(\theta)$  を、次式 (6) の確率密度関数に代入し、確信度  $B_{IPD}(\theta)$  を得る。

【0040】

【数6】

$$B_{IPD}(\theta) = \int_{-\infty}^{X(\theta)} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) dx \quad \dots (6)$$

ここで、 $X(\theta) = (d(\theta) - m) / (s / n)$ 、 $m$  は、 $d(\theta)$  の平均、 $s$  は  $d(\theta)$  の分散であり、 $n$  は IPD の仮説の個数 (本実施形態では 37 個) である。

【0041】

- IID 確信度 -

IID の確信度は、以下のようにして求める。まず、話者  $HM_j$  に対応する調波構造が含む倍音の音圧差の合計を次式 (7) で計算して求める。

【0042】

【数7】

$$S = \sum_{k=0}^{K-1} \Delta \rho(f_k) \quad \dots (7)$$

【0043】

ここで、 $K$  は、その調波構造に含まれる倍音の数を示し、 $\rho(f_k)$  は、IID 計算部 15 で求めた IID である。

次に、表 1 を利用して、音源方向の右らしさ、正面らしさ、左らしさを確信度とする。なお、表 1 は、実験的に得られた値である。

例えば、表 1 を参照して、仮説の音源方位  $\theta$  が  $40^\circ$  で、音圧差  $S$  が正であれば確信度  $B_{IID}(\theta)$  は、左上の欄を参照して 0.35 とする。

【0044】

10

20

30

40

【表 1】

$\theta$		$90^\circ \sim 30^\circ$	$30^\circ \sim -30^\circ$	$-30^\circ \sim -90^\circ$
S	+	0.35	0.5	0.65
	-	0.65	0.5	0.35

## 【0045】

確信度統合部 18

確信度統合部 18 は、Dempster - Shafer 理論に基づき、IPD と IID の確信度  $B_{IPD}(\theta)$ 、 $B_{IID}(\theta)$  を次式 (8) によって統合し、統合確信度  $B_{IPD+IID}(\theta)$  を計算する。そして、統合確信度  $B_{IPD+IID}(\theta)$  が最も大きくなる音源方向  $\theta$  を、話者  $H M_j$  のいる方向とし、以下  $\theta_{H M_j}$  とする。

10

## 【0046】

## 【数 8】

$$B_{IPD+IID}(\theta) = 1 - (1 - B_{IPD}(\theta))(1 - B_{IID}(\theta)) \dots (8)$$

## 【0047】

以上のような聴覚エピポラ幾何を使用した仮説に代えて、頭部伝達関数を用いた仮説データ、又は散乱理論に基づく仮説データを用いることもできる。

20

(頭部伝達関数仮説データ)

頭部伝達関数仮説データは、ロボット周囲から発せられたインパルスより得られる、マイク  $M_R$  とマイク  $M_L$  で検出した音の位相差及び音圧差である。

頭部伝達関数仮説データは、 $-90^\circ$  から  $90^\circ$  の間の適当な間隔 (例えば  $5^\circ$ ) の方向から発したインパルスを、マイク  $M_R$ 、 $M_L$  で検出し、それぞれを周波数分析して周波数  $f$  に対する位相応答及び振幅応答を求め、その差を計算することによって得られる。

得られた頭部伝達関数仮説データは、図 7 (a) の IPD 及び (b) の IID のようになる。

頭部伝達関数を用いる場合には、IPD だけではなく、IID についてもある音源方向から来た音の周波数と IID の関係が求められるので、IPD と IID の両方について距離データ  $d(\theta)$  を作ってから確信度を求める。仮説データの作成方法は、IPD と IID で変わりはない。

30

聴覚エピポラ幾何を利用した仮説データの作成方法と異なり、計算ではなく計測で、各音源方向で発せられた信号に対する周波数  $f$  と IPD の関係を求める。すなわち、図 7 (a)、(b) にある実測値から、それぞれの仮説と入力との距離である  $d(\theta)$  を直接計算する。

## 【0048】

(散乱理論に基づく仮説データ)

散乱理論は、音を散乱する物体、例えばロボットの頭部による散乱波を考慮して、IPD、IID の双方を計算的に推定する理論である。ここでは、音を散乱する物体の内、マイクの入力に主に影響を与える物体はロボットの頭部であると仮定し、これを半径  $a$  の球と仮定する。また頭部の中心の座標を極座標の原点とする。

40

点音源の位置を  $r_0$ 、観測点を  $r$  とすると、観測点における直接音によるポテンシャルは、次式 (9) によって定義される。

## 【数 9】

$$V^i = \frac{v}{2\pi R f} e^{i \frac{2\pi R f}{v}} \dots (9)$$

但し、

50

f : 点音源の周波数

v : 音速

R : 点音源と観測点の距離

また、観測点 r を頭部表面とすると、直接音と散乱音によるポテンシャルは、

J. J. Bowman, T. B. A. Senior, and P. L. E. Uslenghi: *Electromagnetic and Acoustic Scattering by Simple Shapes*. Hemisphere Publishing Co., 1987. などに開示されているように、次式(10)で定義される。

【数10】

$$S(\theta, f) = V^i + V^s \quad \dots (10)$$

$$= - \left( \frac{v}{2\pi a f} \right)^2 \sum_{n=0}^{\infty} (2n+1) P_n(\cos\theta) \frac{h_n^{(1)} \left( \frac{2\pi r_0 f}{v} \right)}{h_n^{(1)'} \left( \frac{2\pi a f}{v} \right)} \dots (10)$$

10

但し、

$V^s$  : 散乱音によるポテンシャル

$P_n$  : 第一種 Legendre 関数

$h_n^{(1)}$  : 第一種球ハンケル関数

$M_R$  の極座標を  $(a, \pi/2, 0)$ 、 $M_L$  の極座標を  $(a, -\pi/2, 0)$  とすると、それぞれにおけるポテンシャルは、次式(11)、(12)で表される。

【数11】

$$S_L(\theta, f) = S\left(\frac{\pi}{2} - \theta, f\right) \quad \dots (11)$$

【数12】

$$S_R(\theta, f) = S\left(-\frac{\pi}{2} - \theta, f\right) \quad \dots (12)$$

20

30

従って、散乱理論に基づく位相差  $IPD_s(\theta, f)$  と音圧差  $IID_s(\theta, f)$  は、それぞれ次式(13)、(14)により求められる。

【数13】

$$\Delta \phi_s(\theta, f) = \arg(S_L(\theta, f)) - \arg(S_R(\theta, f)) \quad \dots (13)$$

【数14】

$$\Delta \rho_s(\theta, f) = 20 \log_{10} \frac{|S_L(\theta, f)|}{|S_R(\theta, f)|} \quad \dots (14)$$

40

【0049】

そして、前記(4)式の  $h_n(\theta, f_k)$  を前記(13)式の  $IPD_s(\theta, f)$  に置き換え、前記した聴覚エピソード幾何を用いた場合と同じ手順で  $B_{IPD}(\theta)$  を求める。

すなわち、 $IPD_s(\theta, f_k)$  と  $IPD_s(\theta, f_k)$  の差を求め、すべてのピーク  $f_k$  について合計して  $d(\theta)$  を求め、得られた  $d(\theta)$  を、前記式(6)の確率密度関数に代入し、確信度  $B_{IPD}(\theta)$  を得る。

【0050】

50

IIDもIPDと同じ方法で $d(\quad)$ と $B_{IID}(\quad)$ を計算する。具体的には、 $h(\quad, f_k)$ を前記(14)式の $IPD_s(\quad, f_k)$ で置き換える。そして、 $s(\quad, f_k)$ と $(f_k)$ の差を求め、すべてのピーク $f_k$ について合計して $d(\quad)$ を求め、得られた $d(\quad)$ を、前記式(6)の確率密度関数に代入し、確信度 $B_{IID}(\quad)$ を得る。

【0051】

このように散乱理論に基づいて音源方向を推定すると、ロボットの頭部の表面に沿って散乱する音声、例えば後頭部を回り込む音の影響を考慮して、音源方向と位相差、および音源方向と音圧差の関係をモデル化できるので、音源方向の推定精度が向上する。特に、音源が側方にある場合は、後頭部を回り込んで音源と反対方向にあるマイクに到達する音のパワーは比較的大きいため、散乱理論を用いることによって音源方向の推定精度が向上する。

10

【0052】

《音源分離部20》

音源分離部20は、音源定位部10により定位された各音源方向の情報、並びに音源定位部で計算したスペクトル(例えばスペクトルCR2)により、各話者HMjの音響(音声)信号を分離する部分である。音源分離方法には、ビームフォーミング、ナルフォーミング、ピーク追跡、指向性マイク、ICA(Independent Component Analysis:独立成分分析)など、従来からある手法を用いることができるが、ここでは、本発明者が開発したアクティブ方向通過型フィルタによる方法について説明する。

20

音源方向の情報を利用して音源を分離する場合、音源の方向がロボットRBの正面から離れるにつれ、2本のマイクを用いて推定した音源方向情報の精度を期待できなくなる。そこで、本実施形態では、正面方向の音源については通過させる方向の範囲を狭く、正面から離れた音源では広くとるように通過帯域をアクティブに制御して、音源の分離精度を向上させる。

【0053】

具体的には、音源分離部20は、図8に示すように、通過帯域関数21と、サブバンド選択部22とを有する。

【0054】

通過帯域関数21

通過帯域関数21は、図9に示したように、音源方向と通過帯域幅の関数で、音源方向が、正面(0°)から離れるにつれ、方向情報の精度を期待できなくなることから、音源方向が正面から離れるほど通過帯域幅が大きくなるように予め設定した関数である。

30

【0055】

サブバンド選択部22

サブバンド選択部22は、スペクトルCR2, CL2の各周波数の値(これを「サブバンド」という)から、特定の方向から来たと推測されるサブバンドを選択する部分である。

サブバンド選択部22では、図10に示すように、音源定位部10で生成した左右の入力音のスペクトルCR2, CL2から、各スペクトルのサブバンドについて、前記式(1)、(2)に従い、 $IPD(f_i)$ 及び $IID(f_i)$ を計算する(図10の両耳間位相差C52, 両耳間音圧差C62参照)。

40

そして、音源定位部10で得られた $H_{Mj}$ を抽出すべき音源方向とし、通過帯域関数21を参照して、 $H_{Mj}$ に対応する通過帯域幅( $H_{Mj}$ )を取得する。取得した通過帯域幅( $H_{Mj}$ )を用いて、通過帯域の最大値 $h$ と最小値 $l$ を次式(15)により求める。通過帯域Bは、方向として平面図で図示すると、例えば図11のようになる。

【0056】

## 【数 15】

$$\left. \begin{aligned} \theta_l &= \theta_{HMj} - \delta(\theta_{HMj}) \\ \theta_h &= \theta_{HMj} + \delta(\theta_{HMj}) \end{aligned} \right\} \dots (15)$$

## 【0057】

次に、 $\theta_l$  と  $\theta_h$  に対応する IPD と IID を推定する。これらの推定には、予め計測、又は計算した伝達関数を利用する。伝達関数は、音源方向から来る信号に対して周波数  $f$  と IPD、IID をそれぞれ関係づけている関数で、前記したエピソード幾何や、頭部伝達関数、散乱理論などを用いる。推定した IPD は、例えば図 10 の両耳間位相差 C 53 における  $\theta_l(f)$ 、 $\theta_h(f)$  であり、推定した IID は、例えば図 10 の両耳間音圧差 C 63 における  $\rho_l(f)$ 、 $\rho_h(f)$  である。

10

## 【0058】

次に、音源方向  $\theta_{HMj}$  に対して、ロボット RB の伝達関数を利用して、スペクトル CR2 または CL2 の周波数  $f_i$  に応じ、周波数  $f_i$  が所定の閾値周波数  $f_{th}$  より小さければ IPD によりサブバンドを選択し、大きければ IID によりサブバンドを選択する。すなわち、以下の条件式 (16) を満たすサブバンドを選択する。

## 【0059】

## 【数 16】

$$\left. \begin{aligned} f_i < f_{th} : \Delta \phi_l(f_i) \leq \Delta \phi(f_i) \leq \Delta \phi_h(f_i) \\ f_i \geq f_{th} : \Delta \rho_l(f_i) \leq \Delta \rho(f_i) \leq \Delta \rho_h(f_i) \end{aligned} \right\} \dots (16)$$

20

## 【0060】

ここで、 $f_{th}$  は、フィルタリングの判断基準に IPD と IID のどちらを用いるかを定める閾値周波数である。

この条件式によれば、例えば、図 10 の両耳間位相差 C 53 においては、周波数  $f_{th}$  より低い周波数で、IPD が  $\theta_l(f)$  と  $\theta_h(f)$  の間にある周波数  $f_i$  のサブバンド (斜線部) が選択される。一方、図 10 の両耳間音圧差 C 63 においては、周波数  $f_{th}$  より高い周波数で、IID が  $\rho_l(f)$  と  $\rho_h(f)$  の間にあるサブバンド (斜線部) が選択される。この選択されたサブバンドからなるスペクトルを本明細書において「選択スペクトル」という。

30

## 【0061】

以上、本実施形態の音源分離部 20 について説明したが、音源分離の方法には、この他に指向性マイクを利用した方法がある。即ち、指向性が狭いマイクをロボット RB に設けておき、音源定位部 10 で得られた音源方向  $\theta_{HMj}$  の方向に指向性マイクを向けるよう、顔の向きを変えれば、その方向から来る音声だけを取得することができる。

この指向性マイクによる方法の場合、1つの指向性マイクしかない場合には、1人の音声しか取得できないという問題もあるが、複数の指向性マイクを所定角度おきに設けておき、音源方向の指向性マイクからの音声信号を利用するようにすれば、複数人の音声の同時取得も可能である。

40

## 【0062】

## 《特徴抽出部 30》

特徴抽出部 30 は、音源分離部 20 で分離された音声スペクトルあるいは分離をしないスペクトル CR2 (または CL2) (以下、音声認識に使用する場合に「認識用スペクトル」という) から音声認識に必要な特徴を抽出する部分である。音声の特徴としては、音声を周波数分析した線形スペクトルや、メル周波数スペクトル、メル周波数ケプストラム係数 (MFCC: Mel-Frequency Cepstrum Coefficient) を用いることができる。本実施形態では、MFCC を用いる場合で説明する。なお、線形スペクトルを特徴として用いる場合は、特徴抽出部 30 は、特に処理を行わない。

50

また、メル周波数スペクトルを用いる場合は、コサイン変換（後述）を行わない。

【 0 0 6 3 】

特徴抽出部 3 0 は、図 1 2 ( a ) に示すように、対数変換部 3 1、メル周波数変換部 3 2、及びコサイン変換部 3 3 を有する。

対数変換部 3 1 は、サブバンド選択部 2 2 ( 図 8 参照 ) が選択した認識用スペクトルの振幅を対数に変換して、対数スペクトルを得る。

メル周波数変換部 3 2 は、対数変換部 3 1 が生成した対数スペクトルを、メル周波数のバンドパスフィルタに通し、周波数がメルスケールに変換されたメル周波数対数スペクトルを得る。

コサイン変換部 3 3 は、メル周波数変換部 3 2 が生成したメル周波数対数スペクトルをコサイン変換する。このコサイン変換により得られた係数が M F C C となる。

10

【 0 0 6 4 】

また、雑音などによって入力音声に変形している場合は、そのスペクトルサブバンドを特徴として信用しないよう、図 1 2 ( b ) に示すように指標 ( 0 から 1 ) を付与するマスキング部 3 4 を、特徴抽出部 3 0 の中または後に任意的に追加してもよい。

図 1 2 ( b ) の例について具体的に説明すると、特徴抽出部 3 0 が任意的にマスキング部 3 4 を含む場合、単語辞書 5 9 は、単語に対応してその単語の時系列スペクトルを有する。ここでは、この時系列スペクトルを「単語音声スペクトル」とする。

単語音声スペクトルは、雑音がない環境下で単語を発声した音声を周波数分析して得られる。特徴抽出部 3 0 に認識用スペクトルが入力されると、入力音声に含まれていると推測された単語の単語音声スペクトルが想定音声スペクトルとして単語辞書から選別される。ここでは、認識用スペクトルと時間長が最も近いものを想定音声スペクトルとして推測する。認識用スペクトルと想定音声スペクトルは、それぞれ対数変換部 3 1、メル周波数変換部 3 2、コサイン変換部 3 3 を経て M F C C に変換される。以下、認識用スペクトルの M F C C を「認識用 M F C C」、想定音声スペクトルの M F C C を「想定 M F C C」とする。

20

マスキング部 3 4 は、認識用 M F C C と想定 M F C C の差を求め、予め想定した閾値より大きい場合は 0 を、小さい場合は 1 を、M F C C の特徴量ベクトルの各特徴ごとに付与する。これを指標として認識用 M F C C と合わせて音声認識部 5 0 に出力する。

想定音声スペクトルを選別する際、1 つだけではなく、複数選別してもよい。また、選別せずに全ての単語音声スペクトルを用いてもよい。その場合には、すべての想定音声スペクトルについて指標を求め、音声認識部 5 0 に出力する。

30

【 0 0 6 5 】

なお、指向性マイクを用いて音源分離を行う場合には、指向性マイクから得られた分離音声に対し、FFT やバンドパスフィルタなどの一般的な周波数分析手法を用いてスペクトルを得る。

【 0 0 6 6 】

《音響モデル合成部 4 0 》

音響モデル合成部 4 0 は、音響モデル記憶部 4 9 に記憶された方向依存音響モデルから、定位された各音源方位に応じた音響モデルを合成する部分である。

40

音響モデル合成部 4 0 は、図 1 3 に示すように、コサイン逆変換部 4 1、線形変換部 4 2、指数変換部 4 3、パラメータ合成部 4 4、対数変換部 4 5、メル周波数変換部 4 6、及びコサイン変換部 4 7 を有し、音響モデル記憶部 4 9 に記憶された方向依存音響モデル  $H(\theta_n)$  を参照して方向の音響モデルを合成する。

【 0 0 6 7 】

音響モデル記憶部 4 9

音響モデル記憶部 4 9 には、ロボット R B の正面を基準とした方向  $\theta_n$  ごとに、方向  $\theta_n$  に適した音響モデルである方向依存音響モデル  $H(\theta_n)$  が記憶されている。方向依存音響モデル  $H(\theta_n)$  は、特定の方向  $\theta_n$  から発せられた人物の音声の特徴を、隠れマルコフモデル ( H M M ) で学習させたものである。各方向依存音響モデル  $H(\theta_n)$  は、図

50

14に示すように、例えば音素を認識単位とし、音素ごとに対応するサブモデル $h(m, n)$ を記憶している。なお、サブモデルは、モノフォン、PTM、バイフォン、トライフォンなど他の認識単位で作成してもよい。

サブモデル $h(m, n)$ の数は、例えば方向 $n$ について $-90^\circ \sim 90^\circ$ まで $30^\circ$ おきに7個のモデルを持ち、サブモデルを40個のモノフォンで構成しているとすれば、合計 $7 \times 40 = 280$ 個となる。

サブモデル $h(m, n)$ は、状態数、各状態の確率密度分布、状態遷移確率の各パラメータを有している。本実施形態では、各音素の状態数は、前部(状態1)、中間部(状態2)、後部(状態3)の3つに固定している。また、本実施形態では、確率密度分布は、正規分布に固定するが、確率密度分布は、正規分布または他の分布の1つ以上の混合分布であってもよい。したがって、本実施形態では、状態遷移確率 $P$ と、正規分布のパラメータ、つまり平均 $\mu$ 及び標準偏差 $\sigma$ を学習させる。

#### 【0068】

サブモデル $h(m, n)$ の学習データは次のようにして作成する。

ロボットRBに対し、音響モデルを作成したい方向から、特定の音素からなる音声信号を図示しないスピーカにより発する。そして、検出した音響信号を特徴抽出部30によりMFCCに変換し、後述する音声認識部50で音声認識させる。すると、認識した音声、音素ごとにどのくらいの確率であるかが結果として得られるが、この結果に対し、特定の方向の特定の音素であるという教師信号を与えることで音響モデルを適応学習させる。そして、サブモデルを学習するのに十分な種類(例えば、異なる話者)の音素や単語を学習させる。

なお、学習用音声を発する際、音響モデルを作成したい方向とは異なる方向から、別の音声をノイズとして発してもよい。この場合は、前記した音源分離部20により音響モデルを作成したい方向の音響のみを分離した上で、特徴抽出部30によりMFCCに変換する。また、これらの学習は、音響モデルを不特定話者のモデルとして持たせたい場合には、不特定の話者の声で学習させればよいし、特定話者ごとにモデルを持たせたい場合には、特定話者ごとに学習させればよい。

#### 【0069】

コサイン逆変換部41から指数変換部43は、確率密度分布のMFCCを線形スペクトルに戻す。つまり、確率密度分布について、特徴抽出部30と逆の操作をする。

#### 【0070】

コサイン逆変換部41

コサイン逆変換部41は、音響モデル記憶部49が記憶している方向依存音響モデル $H(n)$ が有するMFCCについてコサイン逆変換してメル対数スペクトルを生成する。

#### 【0071】

線形変換部42

線形変換部42は、コサイン逆変換部41により生成されたメル対数スペクトルの周波数を線形周波数に変換し、対数スペクトルを生成する。

#### 【0072】

指数変換部43

指数変換部43は、線形変換部42により生成された対数スペクトルの強度を指数変換し、線形スペクトルを生成する。線形スペクトルは、平均 $\mu$ 、標準偏差 $\sigma$ の確率密度分布として得られる。

#### 【0073】

パラメータ合成部44

パラメータ合成部44は、図15に示すように、方向依存音響モデル $H(n)$ にそれぞれ重みをかけた上でそれらの和をとり、音源方向 $H_{mj}$ の音響モデル $H(H_{mj})$ を合成する。方向依存音響モデル $H(n)$ にある各サブモデルは、それぞれコサイン逆変換部41から指数変換部43により、線形スペクトルの確率密度分布に変換され、それぞれ、平均 $\mu_{1nm}, \mu_{2nm}, \mu_{3nm}$ 、標準偏差 $\sigma_{1nm}, \sigma_{2nm}, \sigma_{3nm}$ 、状

10

20

30

40

50

状態遷移確率  $P_{11nm}, P_{12nm}, P_{22nm}, P_{23nm}, P_{33nm}$  のパラメータを持っている。そして、これらのパラメータを、予め学習によって求められ、音響モデル記憶部 49 に記憶されている重みと内積して、音源方向  $\theta_{HMj}$  の音響モデルを合成する。つまり、パラメータ合成部 44 は、方向依存音響モデル  $H(\theta_n)$  の線形和により音源方向  $\theta_{HMj}$  の音響モデルを合成している。なお、重み  $W_{n\theta_{HMj}}$  の設定の仕方は後述する。

【0074】

$H(\theta_{HMj})$  にあるサブモデルを合成する場合には、状態 1 の平均  $\mu_{1\theta_{HMj}m}$  を次式 (17) により求める。

【0075】

【数17】

$$\mu_{1\theta_{HMj}m} = \frac{1}{\sum_{n=1}^N W_{n\theta_{HMj}}} \sum_{n=1}^N W_{n\theta_{HMj}} \mu_{1nm} \dots (17)$$

10

【0076】

平均  $\mu_{2\theta_{HMj}m}, \mu_{3\theta_{HMj}m}$  についても同様にして求めることができる。

【0077】

また、状態 1 の標準偏差  $\sigma_{1\theta_{HMj}m}$  の合成については、共分散  $\sigma_{1\theta_{HMj}m}^2$  を次式 (18) により求める。

20

【数18】

$$\sigma_{1\theta_{HMj}m}^2 = \frac{1}{\sum_{n=1}^N W_{n\theta_{HMj}}} \sum_{n=1}^N W_{n\theta_{HMj}} \sigma_{1nm}^2 \dots (18)$$

【0078】

標準偏差  $\sigma_{2\theta_{HMj}m}, \sigma_{3\theta_{HMj}m}$  についても同様にして求めることができる。得られた  $\mu$  と  $\sigma$  により、確率密度分布を求めることができる。

30

【0079】

また、状態 1 の状態遷移確率  $P_{11\theta_{HMj}m}$  の合成については、次式 (19) により求める。

【0080】

【数19】

$$P_{11\theta_{HMj}m} = \frac{1}{\sum_{n=1}^N W_{n\theta_{HMj}}} \sum_{n=1}^N W_{n\theta_{HMj}} P_{11nm} \dots (19)$$

40

【0081】

状態遷移確率  $P_{12\theta_{HMj}m}, P_{22\theta_{HMj}m}, P_{23\theta_{HMj}m}, P_{33\theta_{HMj}m}$  についても同様にして求めることができる。

【0082】

次に、対数変換部 45 からコサイン変換部 47 により、確率密度分布を線形スペクトルから MFCC に変換し直す。すなわち、対数変換部 45 は、対数変換部 31 と、メル周波数変換部 46 は、メル周波数変換部 32 と、コサイン変換部 47 は、コサイン変換部 33 と同様であるので、詳細な説明を省略する。

【0083】

なお、単一正規分布ではなく、混合正規分布の形で合成する場合には、前記した平均  $\mu$

50



、標準偏差 の計算に代えて次式 ( 2 0 ) により確率密度分布  $f_{1 \theta_{HMj} m} ( x )$  を求める。

【 0 0 8 4 】

【 数 2 0 】

$$f_{1 \theta_{HMj} m} ( x ) = \frac{1}{\sum_{n=1}^N W_{n \theta_{HMj}}} \sum_{n=1}^N W_{n \theta_{HMj}} f_{1 n m} ( x ) \dots ( 2 0 )$$

【 0 0 8 5 】

10

確率密度分布  $f_{2 \theta_{HMj} m} ( x )$  ,  $f_{3 \theta_{HMj} m} ( x )$  についても同様にして求めることができる。

【 0 0 8 6 】

パラメータ合成部 4 4 は、このようにして得られた音響モデルを、音響モデル記憶部 4 9 に記憶させる。

なお、このような音響モデルの合成は、音声認識装置 1 が作動している間、パラメータ合成部 4 4 がリアルタイムに行う。

【 0 0 8 7 】

重み  $W_{n \theta_{HMj}}$  の設定

重み  $W_{n \theta_{HMj}}$  は、音源方向  $\theta_{HMj}$  に対応する音響モデルを合成するとき、各方向依存音響モデル  $H ( \theta_n )$  に対して設定するもので、 $H ( \theta_n )$  に含まれるすべてのサブモデル  $h ( m, \theta_n )$  に対して用いる重み  $W_{n \theta_{HMj}}$  を設定してもよいし、あるいは各サブモデル  $h ( m, \theta_n )$  に対応する重み  $W_{m n \theta_{HMj}}$  を設定してもよい。基本的には、音源が正面にある場合の重み  $W_{n 0}$  を定める関数  $f ( \theta )$  をあらかじめ設定しておく、音源方向  $\theta_{HMj}$  に対応する音響モデルを合成する際に、 $f ( \theta )$  を  $\theta$  軸方向に  $\theta_{HMj}$  移動 (  $\theta - \theta_{HMj}$  とする ) した関数  $f ( \theta - \theta_{HMj} )$  を求め、これを参照して  $W_{n \theta_{HMj}}$  を設定する。

20

【 0 0 8 8 】

関数  $f ( \theta )$  の作成

[ A ]  $f ( \theta )$  を経験的に求める方法

30

$f ( \theta )$  を経験的に求める場合は、経験的に得られた定数  $a$  を用いて次式のように表す。

$$f ( \theta ) = a + \theta \quad ( \theta < 0 , \theta = - 9 0 ^\circ \text{ のとき } f ( \theta ) = 0 )$$

$$f ( \theta ) = - a + \theta \quad ( \theta \geq 0 , \theta = 9 0 ^\circ \text{ のとき } f ( \theta ) = 0 )$$

ここで、定数  $a = 1 . 0$  とすれば、音源が正面にある場合の  $f ( \theta )$  は、図 1 6 ( a ) のようになる。また、 $f ( \theta )$  を  $\theta$  軸方向に  $\theta_{HMj}$  移動したのが図 1 6 ( b ) である。

【 0 0 8 9 】

[ B ]  $f ( \theta )$  を学習によって求める方法

$f ( \theta )$  を学習によって求める場合は、例えば次のような学習をする。

音源が正面にあるときの任意の音素  $m$  の重みを  $W_{m n 0}$  とする。最初に適当な初期値の重みの値の  $W_{m n 0}$  を設定しておく、この  $W_{m n 0}$  を用いて合成した音響モデル  $H ( \theta_0 )$  で  $m$  を含む適当な音素列、例えば音素列 [  $m m m$  ] を認識させる試行を行う。具体的には、正面に設置したスピーカから、前記音素列を発生し、これを認識させる。ここで、学習データは、1つの音素  $m$  自体であってもよいのであるが、音素が複数つながった音素列で学習させた方がよい学習結果が得られるため、音素列を使用している。

40

この時の認識結果が、例えば図 1 7 である。図 1 7 では、初期値の  $W_{m n 0}$  を用いて合成した音響モデル  $H ( \theta_0 )$  での認識結果が 1 行目であり、2 行目以下の  $H ( \theta_n )$  が方向  $\theta_n$  の方向依存音響モデル  $H ( \theta_n )$  を使用したときの認識結果である。例えば、音響モデル  $H ( \theta_{90} )$  での認識結果は音素列 [ / x / / y / / z / ] であり、音響モデル  $H ( \theta_0 )$  での認識結果は、音素列 [ / x / / y / m ] であったことを示す。

50

1回目の試行後、まず1音素目を見て、図17の正面から  $\theta = \pm 90^\circ$  の範囲に一致する音素が認識された場合、その方向に対応するモデルの重み  $W_{m n \theta_0}$  を  $d$  増加させる。 $d$  は実験的に求め、例えば0.05とする。そして、一致する音素が認識されない場合、その方向に対応するモデルの重み  $W_{m n \theta_0}$  を  $d / (n - k)$  減少させる。つまり、正解を出した方向依存音響モデルの重みは大きくし、正解を出さなかった方向依存音響モデルの重みは減少させる。

【0090】

例えば、図17の場合では、 $H(n)$  と  $H(\theta_0)$  が一致しているので、対応する重み  $W_{m n}$  と重み  $W_{m \theta_0}$  を  $d$  増加させ、それ以外の重みを  $2d / (n - 2)$  減少させる。

10

一方、1音素目に一致する音素を認識した方向  $n$  が無い場合、他の方向に対して重みの大きい、優勢な方向依存音響モデル  $H(n)$  があれば、その方向依存音響モデル  $H(n)$  の重みを  $d$  減少させ、それ以外のモデルの重みを  $k d / (n - k)$  増加させる。つまり、どの方向依存音響モデル  $H(n)$  も認識できなかったということは、現在の重みの分配が良くない可能性があるから、現在の重みが優勢な方向について重みを減少させる。

優勢であるかどうかは、重みが予め定められた閾値（ここでは0.8とする）より大きいかどうかで判断する。優勢な方向依存音響モデル  $H(n)$  がなければ、最大の重みのみを  $d$  減少させ、その他の方向依存音響モデル  $H(n)$  の重みを  $d / (n - 1)$  増加させる。

20

そして、更新された重みを用いて、前記した試行を繰り返す。

そして、音響モデル  $H(\theta_0)$  の認識結果が、正解  $m$  となったときに、繰り返しを終了し、次の音素  $m$  の認識および学習へ移るか、または学習を終了する場合、ここで得られた重み  $W_{m n \theta_0}$  が  $f(\quad)$  となる。次の音素  $m$  へ移る場合は、すべての音素について学習し、得られた  $W_{m n \theta_0}$  を平均したものが  $f(\quad)$  となる。

これを平均せず、各サブモデル  $h(m, n)$  に対応する重み  $W_{m n H M j}$  を  $f(\quad)$  にしてもよい。

なお、所定の回数（例えば  $0.5 / d$  回）繰り返しても、音響モデル  $H(H M j)$  の認識結果が正解に至らない場合、例えば  $m$  の認識がうまくいかなかった場合には、次の音素  $m$  の学習へ移り、最終的にうまく認識できた音素（例えば  $m$ ）の重みの分布と同じ値で重みを更新する。

30

また、音響モデルを合成するたびに  $f(\quad - H M j)$  を求めるのではなく、予め適当な  $H M j$  について、 $H(n)$  に含まれるすべてのサブモデル  $h(m, \quad)$ （表2参照）が用いる重み  $W_{n H M j}$  または各サブモデル  $h(m, n)$  に対応する  $W_{n H M j}$  を求めた表3を作成しておいてもよい。なお、表2および表3において、添え字の  $1 \cdot \cdot \cdot m \cdot \cdot \cdot M$  は音素を表し、 $1 \cdot \cdot \cdot n \cdot \cdot \cdot N$  は方向を表す。

【表 2】

$H(\theta_1)$	$H(\theta_2)$	...	$H(\theta_n)$	...	$H(\theta_N)$
$h(1, \theta_1)$	$h(1, \theta_2)$	...	$h(1, \theta_n)$	...	$h(1, \theta_N)$
$\vdots$	$\vdots$	...	$\vdots$	...	$\vdots$
$h(m, \theta_1)$	$h(m, \theta_2)$	...	$h(m, \theta_n)$	...	$h(m, \theta_N)$
$\vdots$	$\vdots$	...	$\vdots$	...	$\vdots$
$h(M, \theta_1)$	$h(M, \theta_2)$	...	$h(M, \theta_n)$	...	$h(M, \theta_N)$

10

【表 3】

$W_1$	$W_2$	...	$W_n$	...	$W_N$
$w_{11}$	$w_{12}$	...	$w_{1n}$	...	$w_{1N}$
$\vdots$	$\vdots$	...	$\vdots$	...	$\vdots$
$w_{m1}$	$w_{m2}$	...	$w_{mn}$	...	$w_{mN}$
$\vdots$	$\vdots$	...	$\vdots$	...	$\vdots$
$w_{M1}$	$w_{M2}$	...	$w_{Mn}$	...	$w_{MN}$

20

30

【0091】

このようにして学習して得られた重みは、音響モデル記憶部49に記憶させる。

【0092】

《音声認識部50》

音声認識部50は、音源方向  $H_{Mj}$  に対応して合成された音響モデル  $H(H_{Mj})$  を用いて、分離された各話者  $H_{Mj}$  の音声あるいは入力音声から抽出した特徴を認識して文字情報とし、単語辞書59を参照して言葉を認識し、認識結果を出力する。この音声認識の方法は一般的な隠れマルコフモデルを利用した認識方法なので、詳細な説明は省略する。

40

なお、マスキング部を特徴抽出部30の中または後に設けて、MFCCの各サブバンドの信用度を示す指標  $\omega(i)$  が付与されている場合には、音声認識部50は、入力された特徴に次式(21)のような処理を行ってから認識する。

【数21】

$$\left. \begin{aligned} x_r &= 1 - x_n \\ x_n(i) &= x(i) \times \omega(i) \end{aligned} \right\} \dots (16)$$

$x_r$  : 音声認識に用いる特徴  
 $x$  : MFCC

50

$i$  : MFCCの成分  
 $x_n$  :  $x$ のうち信用できない成分

そして、得られた出力確率と状態遷移確率を用いて、一般的な隠れマルコフモデルを利用した認識方法と同様に認識を行う。

【0093】

以上のように構成された、音声認識装置1による動作を説明する。

図1に示すように、ロボットRBのマイク $M_R, M_L$ に、複数の話者 $HM_j$  (図3参照)の音声が入力される。

そして、マイク $M_R, M_L$ が検出した音響信号の音源方向が音源定位部10で定位される。音源定位は、前記したように周波数分析、ピーク抽出、調波構造の抽出、IPD・IIDの計算の後、聴覚エピソード幾何に基づいた仮説データを利用して確信度を計算する。そして、IPDとIIDの確信度を統合して最も可能性が高い $HM_j$ を音源方向とする(図2参照)。

10

【0094】

次に、音源分離部20で、音源方向 $HM_j$ の音を分離する。音源分離は、通過帯域関数を利用して、音源方向 $HM_j$ のIPD及びIIDのそれぞれの上限值 $h(f)$ ,  $l(f)$ 及び下限値 $l(f)$ ,  $l(f)$ を求め、前記式(16)の条件と、この上限値、下限値の条件とから、音源方向 $HM_j$ のスペクトルと推定されるサブバンド(選択スペクトル)を選択する。その後、選択サブバンドのスペクトルを逆FFTにより変換すれば、音声信号に変換できる。

20

【0095】

次に、特徴抽出部30は、音源分離部20が分離した選択スペクトルを、対数変換部31、メル周波数変換部32、コサイン変換部33によりMFCCに変換する。

【0096】

一方、音響モデル合成部40は、音響モデル記憶部49に記憶された方向依存音響モデル $H(\theta_n)$ と、音源定位部10が定位した音源方向 $HM_j$ とから、音源方向 $HM_j$ に適切と考えられる音響モデルを合成する。

すなわち、音響モデル合成部40は、方向依存音響モデル $H(\theta_n)$ を、コサイン逆変換部41、線形変換部42、及び指数変換部43により、線形スペクトルに変換する。そして、パラメータ合成部44は、音源方向 $HM_j$ の重み $W_{n, HM_j}$ を音響モデル記憶部49から読み出し、これと方向依存音響モデル $H(\theta_n)$ との内積をとって、音源方向 $HM_j$ の音響モデル $H(\theta_{HM_j})$ を合成する。そして、この線形スペクトルで表された音響モデル $H(\theta_{HM_j})$ を、対数変換部45、メル周波数変換部46、及びコサイン変換部47によりMFCCで表した音響モデル $H(\theta_{HM_j})$ に変換する。

30

【0097】

次に、音声認識部50は、音響モデル合成部40で合成された音響モデル $H(\theta_{HM_j})$ を利用して、隠れマルコフモデルにより音声認識を行う。

【0098】

このようにして、音声認識を行った結果の例が、表4である。

【0099】

【表4】

40

音響モデルの方向	従来手法							本発明
	-90°	-60°	-30°	0	30°	60°	90°	
孤立単語認識率	20%	20%	38%	42%	60%	59%	50%	78%

【0100】

表4に示すように、方向依存音響モデルを-90°~90°まで30°おきに用意して、各音響モデルで40°の方向から孤立単語を認識させたところ(従来手法)、最も認識率が高くて30°方向の方向依存音響モデルを用いた60%であった。これに対し、本実施形態の手法を使用して40°方向の音響モデルを合成して、これを用いて孤立単語を

50

認識させたところ、78%の高い認識率を示した。このように、本実施形態の音声認識装置1によれば、任意の方向から音声が発せられた場合であっても、その方向に適した音響モデルをその都度合成するので、高い認識率を実現することができる。また、任意の方向の音声を認識できることから、移動している音源からの音声認識や、移動体(ロボットRB)自身が移動しているときにも、高い認識率での音声認識が可能である。

【0101】

また、方向依存音響モデルを、断続的な数個、例えば音源方向にして60°ごとや30°ごとに記憶しておけば良く、音響モデルの学習に必要なコストを小さくすることができる。

さらに、合成した音響モデル一つについて音声認識を行えば良いため、複数方向の音響モデルについて音声認識を試みる並列処理も不要であり、計算コストを小さくすることができる。そのため、実時間処理や、組み込み用途には好適である。

10

【0102】

以上、本発明の第1実施形態について説明したが、本発明は第1実施形態には限定されず、以下の実施形態のように変形して実施することが可能である。

【0103】

[第2実施形態]

第2実施形態では、第1実施形態の音源定位部10に代えて、相互相関のピークを用いて音源方向を定位する音源定位部110を備える。なお、他の部分については第1実施形態と同様であるので説明を省略する。

20

《音源定位部110》

第2実施形態に係る音源定位部110は、図18に示すように、フレーム切り出し部111、相互相関計算部112、ピーク抽出部113、方向推定部114を有する。

【0104】

フレーム切り出し部111

フレーム切り出し部111は、左右のマイクM<sub>R</sub>, M<sub>L</sub>に入力されたそれぞれの音響信号について、所定の時間長、例えば100msで切り出す処理を行う。切り出し処理は、適当な時間間隔、例えば30msごとに行われる。

【0105】

相互相関計算部112

30

相互相関計算部112は、フレーム切り出し部111が切り出した右マイクM<sub>R</sub>の音響信号と、左マイクM<sub>L</sub>の音響信号とで、次式(22)により相互相関を計算する

【数22】

$$CC(T) = \int_0^T x_L(t)x_R(t+T)dt \quad \dots (22)$$

但し、

CC(T) : x<sub>L</sub>(t)とx<sub>R</sub>(t)の相互相関

T : フレーム長

x<sub>L</sub>(t) : フレーム長Tで切り出された、マイクLからの入力信号

40

x<sub>R</sub>(t) : フレーム長Tで切り出された、マイクRからの入力信号

【0106】

ピーク抽出部113

ピーク抽出部113は、得られた相互相関の結果からピークを抽出する。抽出するピークの数、音源の数が予め分かっている場合は、その数に対応したピークを大きいものから選択する。音源数が不明なときは、予め定めた閾値を超えたピークを全て抽出するか、あるいは予め定めた所定数のピークを大きいものから順に選択する。

【0107】

方向推定部114

音源方向 H<sub>Mj</sub> は、得られたピークから、右マイクM<sub>R</sub>と左マイクM<sub>L</sub>に入力された

50

音響信号の到達時間差  $D$  に音速  $v$  を掛けて、図 19 に示す距離差  $d$  を計算し、さらに、次式により求める。

$$\theta_{HMj} = \arcsin(d / 2r)$$

【0108】

このような相互相関を用いた音源定位部 110 によっても、音源方向  $\theta_{HMj}$  の方向が推定され、前記した音響モデル合成部 40 により、音源方向  $\theta_{HMj}$  に適した音響モデルを合成することで、認識率の向上を図ることができる。

【0109】

[第3実施形態]

第3実施形態では、第1実施形態に加えて、音源定位部音源が同一音源から来ていることを確認しながら音声認識を行う機能を追加している。なお、第1実施形態と同じ部分については、同じ符号を付して説明を省略する。

第3実施形態に係る音声認識装置 100 は、図 20 に示すように、第1実施形態の音声認識装置 1 に加え、音源定位部 10 が定位した音源方向を入力されて、音源を追跡し、同じ音源から音響が来続けているかを確認し、確認ができたなら、音源方向を音源分離部 20 へ出力するストリーム追跡部 60 を有している。

【0110】

図 21 に示すように、ストリーム追跡部 60 は、音源方向履歴記憶部 61 と、予測部 62 と、比較部 63 とを有する。

【0111】

音源方向履歴記憶部 61 は、図 22 に示すような、時間と、その時間において認識された音源の方向及び音源のピッチ（その音源の調波構造が持つ基本周波数  $f_0$ ）とが関連付けて記憶されている。

【0112】

予測部 62 は、音源方向履歴記憶部 61 から、直前まで追跡していた音源の音源方向の履歴を読み出し、直前までの履歴からカルマンフィルタなどにより現時点  $t_1$  での音源方向  $\theta_{HMj}$  及び基本周波数  $f_0$  とからなるストリーム特徴ベクトル  $(\theta_{HMj}, f_0)$  を予測し、比較部 63 へ出力する。

【0113】

比較部 63 は、音源定位部 10 から、音源定位部 10 で定位された現時点  $t_1$  の各話者  $HMj$  の音源方向  $\theta_{HMj}$  と、その音源の基本周波数  $f_0$  とが入力される。そして、予測部 62 から入力された予測したストリーム特徴ベクトル  $(\theta_{HMj}, f_0)$  と、音源定位部 10 で定位された音源方向及びピッチから求まるストリーム特徴ベクトル  $(\theta_{HMj}, f_0)$  を比較して、その差（距離）が予め定めた閾値よりも小さい場合に、音源方向  $\theta_{HMj}$  を音源分離部へ出力する。また、ストリーム特徴ベクトル  $(\theta_{HMj}, f_0)$  を音源方向履歴記憶部 61 へ記憶させる。

前記した差（距離）が、予め定めた閾値よりも大きい場合には、定位した音源方向  $\theta_{HMj}$  を音源分離部 20 へ出力しないので、音声認識は行われない。なお、音源方向  $\theta_{HMj}$  とは別に、音源の追跡ができていないか否かを示すデータを、比較部 63 から音源分離部 20 へ出力してもよい。

なお、基本周波数  $f_0$  を用いず、音源方向  $\theta_{HMj}$  だけで予測してもよい。

【0114】

このようなストリーム追跡部 60 を有する音声認識装置 100 によれば、音源定位部 10 で音源方向が定位され、ストリーム追跡部 60 へ音源方向とピッチが入力される。ストリーム追跡部 60 では、予測部 62 が、音源方向履歴記憶部 61 に記憶された音源方向の履歴を読み出して現時点  $t_1$  でのストリーム特徴ベクトル  $(\theta_{HMj}, f_0)$  を予測する。比較部 63 は、予測部 62 で予測されたストリーム特徴ベクトル  $(\theta_{HMj}, f_0)$  と、音源定位部 10 から入力された値から求まるストリーム特徴ベクトル  $(\theta_{HMj}, f_0)$  とを比較して、その差（距離）が所定の閾値より小さければ、音源方向を音源分離部 20 へ出力する。

10

20

30

40

50

音源分離部 20 は、音源定位部 10 から入力されたスペクトルのデータと、ストリーム追跡部 60 が出力した音源方向  $H_{M_j}$  のデータに基づき、第 1 実施形態と同様にして音源を分離する。そして、以下、特徴抽出部 30、音響モデル合成部 40、音声認識部 50 でも、第 1 実施形態と同様にして、処理を行う。

【0115】

このように、本実施形態の音声認識装置 100 は、音源が追跡できているか否かを確認した上で音声認識を行うので、音源が移動している場合にも、同じ音源が発し続けている音声を連続して認識するため、誤認識の可能性を低くすることができる。特に、複数の移動する音源があって、それらの音源が交差する場合などに好適である。

また、音源方向を記憶、予測していることから、その方向の所定範囲についてのみ音源を探索すれば、処理を少なくすることができる。

【0116】

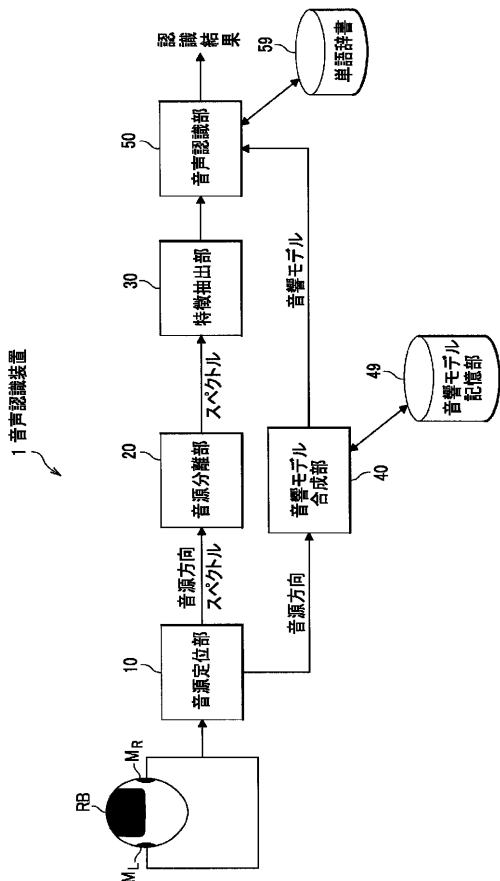
以上、本発明の実施形態について説明したが、本発明は、前記した実施形態には限定されず適宜変更して実施される。

例えば、音声認識装置 1 が、カメラと、公知の画像認識装置を有し、話者の顔を認識して、誰が話しているかを自己が有するデータベースから話者を特定する話者同定部を備え、前記方向依存音響モデルを話者ごとに有していれば、話者に適した音響モデルを合成することができるので、認識率をより高くする事ができる。あるいは、カメラを使わず、ベクトル量子化 (VQ) を用いて、予め登録してある話者の音声をベクトル化したものと、音源分離部 20 で分離された音声をベクトル化したものとを比較し、最も距離の近い話者を結果として出力することで話者を同定してもよい。

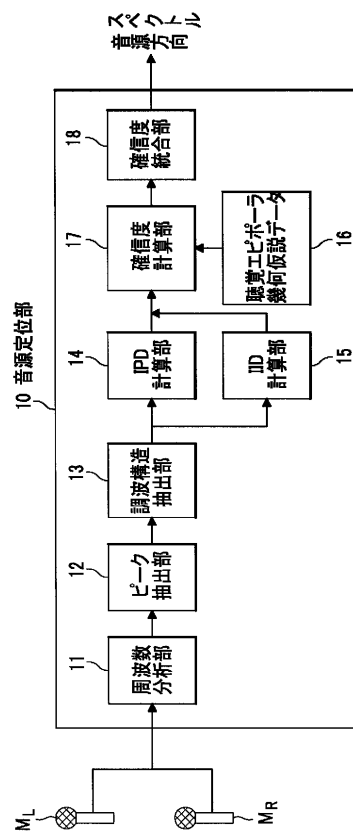
10

20

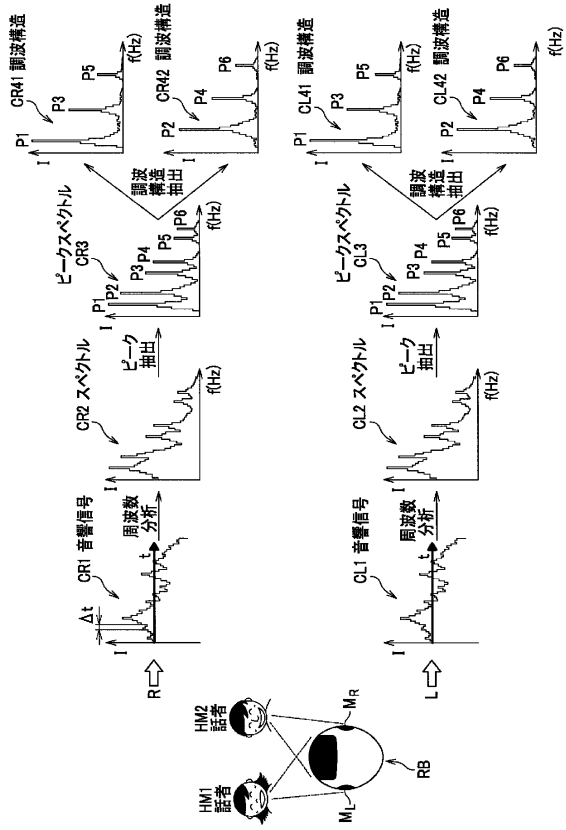
【図 1】



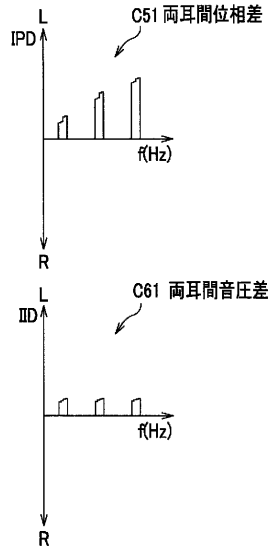
【図 2】



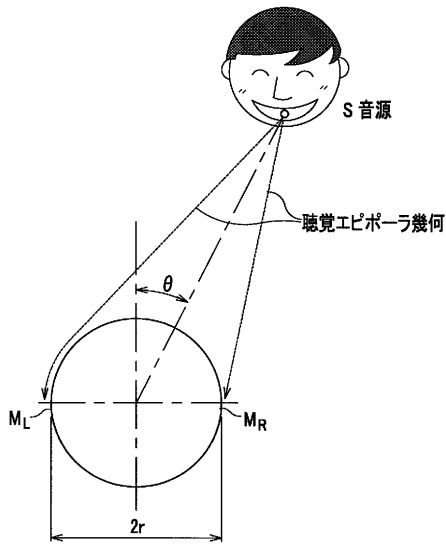
【 図 3 】



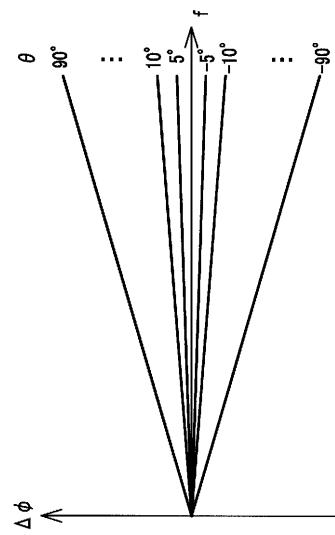
【 図 4 】



【 図 5 】

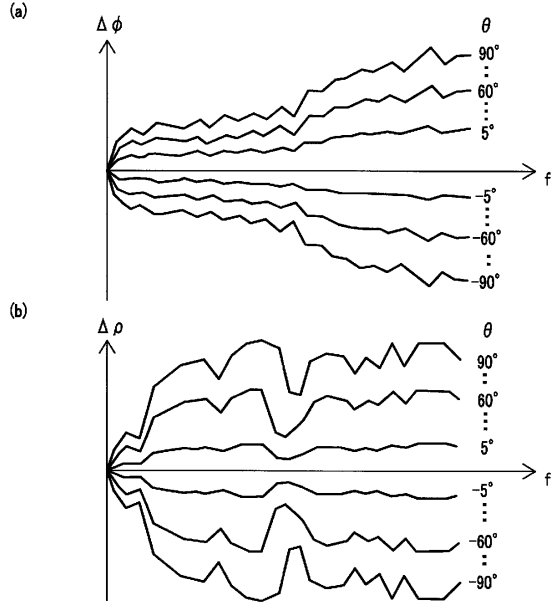


【 図 6 】

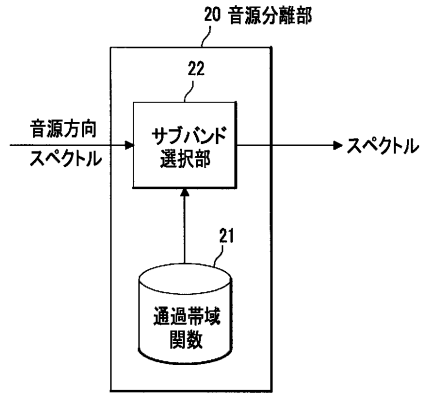




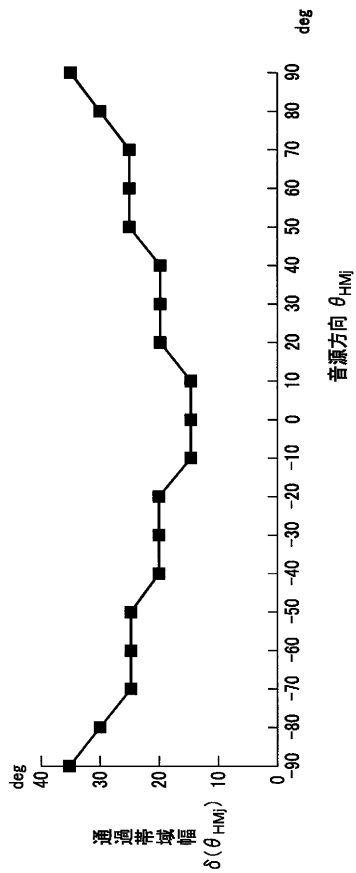
【 図 7 】



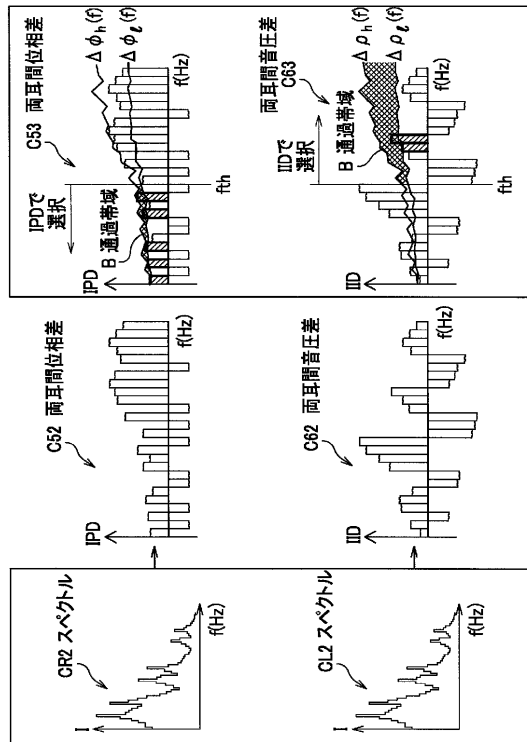
【 図 8 】



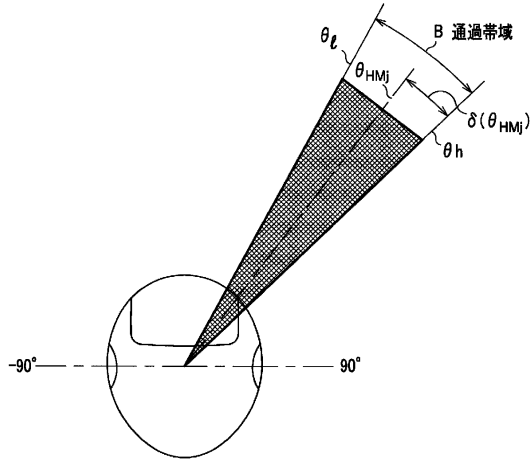
【 図 9 】



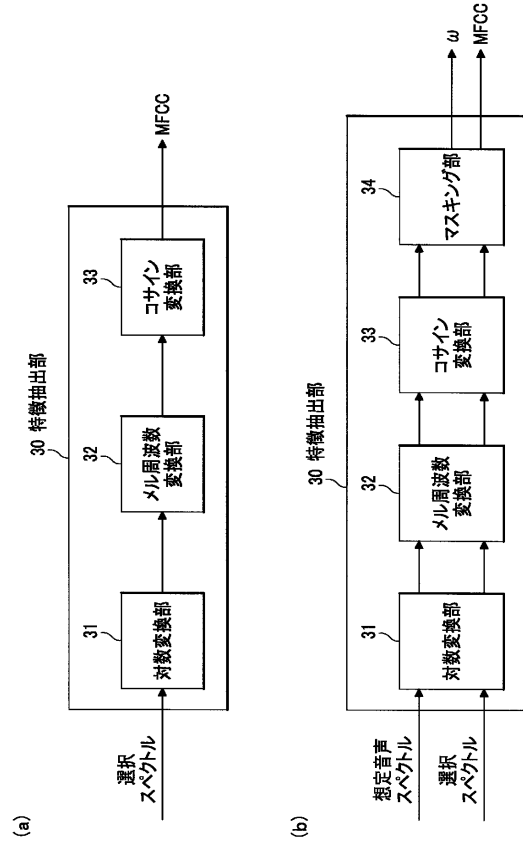
【 図 10 】



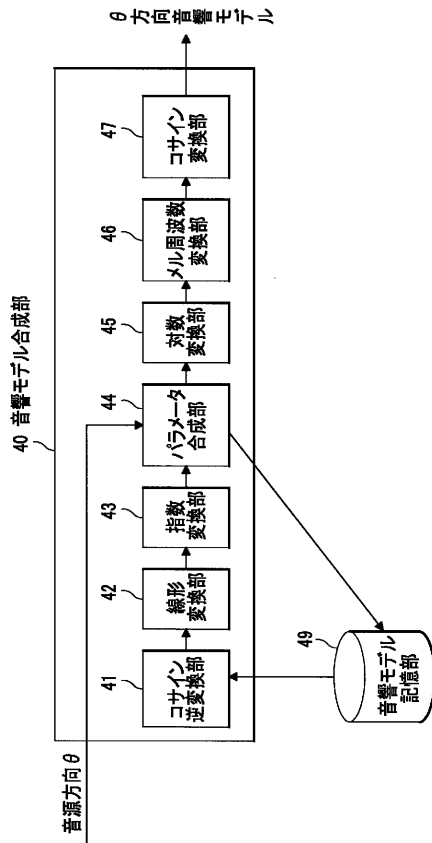
【図 1 1】



【図 1 2】



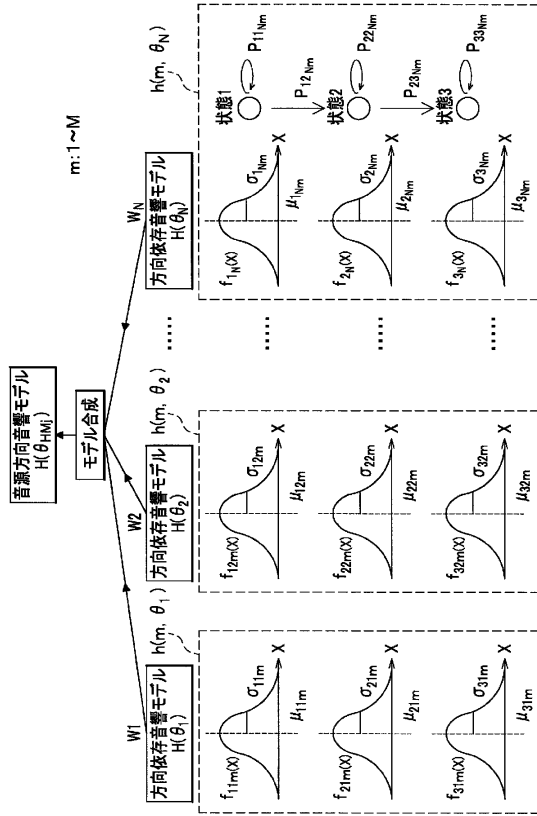
【図 1 3】



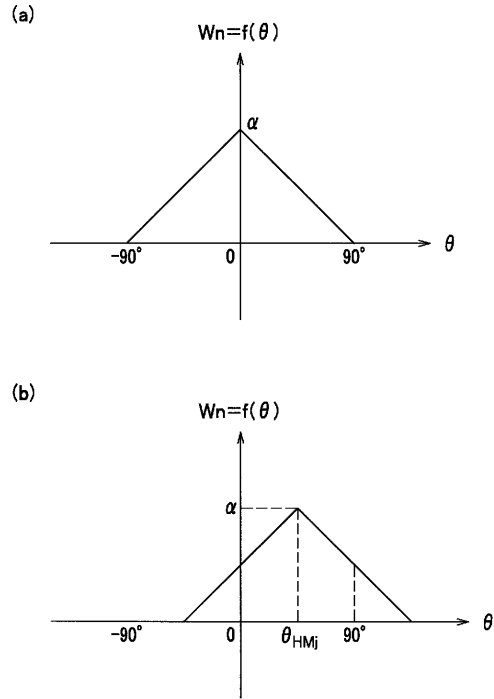
【図 1 4】

認識単位	サブモデル
/a/	$h(/a/, \theta_n)$
/b/	$h(/b/, \theta_n)$
⋮	⋮
m	$h(m, \theta_n)$
⋮	⋮

【図15】



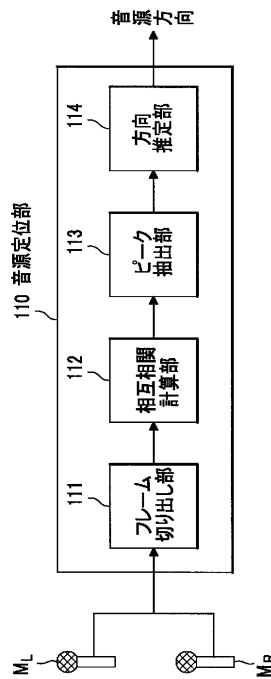
【図16】



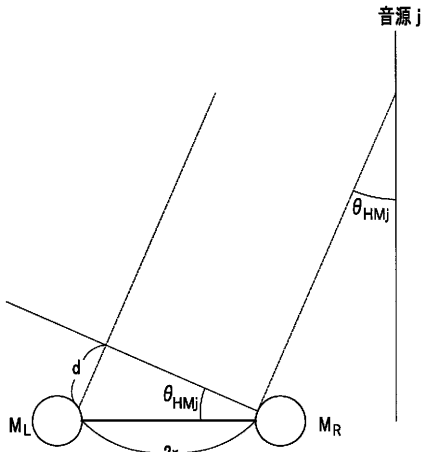
【図17】

認識単位	m	m'	m''
$H(\theta_{HMj})$	/x/	/y/	/z/
$H(\theta_{-90})$	/x/	/y/	m''
⋮	⋮	⋮	⋮
$H(\theta_n)$	m	/y/	/z/
⋮	⋮	⋮	⋮
$H(\theta_{90})$	m	/y/	m''

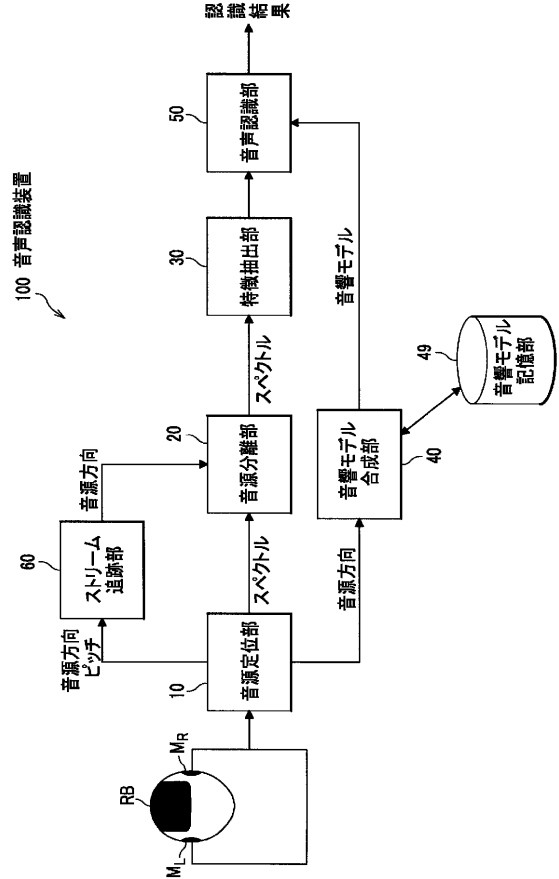
【図18】



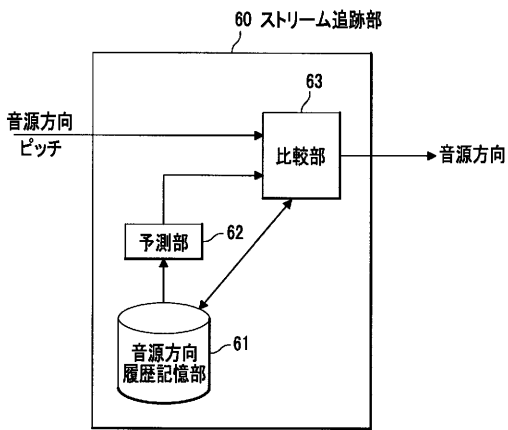
【図19】



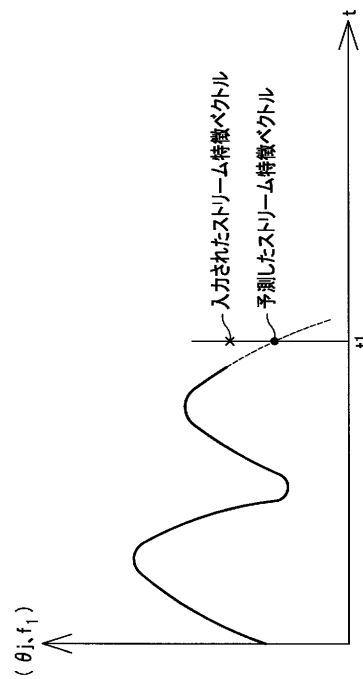
【図20】



【図21】



【図22】



---

フロントページの続き

審査官 菊池 智紀

- (56)参考文献 特開平11-143486(JP,A)  
特開2000-066698(JP,A)  
特開2002-264051(JP,A)  
特表2001-511267(JP,A)  
特開2002-041079(JP,A)  
特開2003-337594(JP,A)  
中臺一博 他, "階層的な視聴覚統合と散乱理論を利用したロボットによる三話者同時発話認識の向上", 日本ロボット学会学術講演会予稿集(CD-ROM), 2003年 9月20日, Vol.21, p.2K14  
中臺一博 他, "アクティブオーディションによる複数音源の定位・分離・認識", 人工知能学会AIチャレンジ研究会資料, 2002年11月22日, Vol.16th, p.25-32  
Kazuhiro NAKADAI et al., "Robot Recognizes Three Simultaneous Speech By Active Audition", Proc. of the 2003 IEEE, 2003年 9月14日, Vol.1, p.398-405

(58)調査した分野(Int.Cl., DB名)

G10L 15/00-15/28

IEEE Xplore

JSTPlus(JDreamII)