

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4157581号
(P4157581)

(45) 発行日 平成20年10月1日(2008.10.1)

(24) 登録日 平成20年7月18日(2008.7.18)

(51) Int. Cl.	F I		
G 1 0 L 15/28 (2006.01)	G 1 0 L 15/28	4 0 0	
G 1 0 L 21/02 (2006.01)	G 1 0 L 21/02	2 0 2 A	
H 0 4 R 3/00 (2006.01)	G 1 0 L 21/02	2 0 1 C	
	H 0 4 R 3/00	3 2 0	

請求項の数 4 (全 17 頁)

(21) 出願番号	特願2006-546764 (P2006-546764)	(73) 特許権者	000005326
(86) (22) 出願日	平成17年12月2日(2005.12.2)		本田技研工業株式会社
(86) 国際出願番号	PCT/JP2005/022601		東京都港区南青山二丁目1番1号
(87) 国際公開番号	W02006/059806	(74) 代理人	110000246
(87) 国際公開日	平成18年6月8日(2006.6.8)		特許業務法人オカダ・フシミ・ヒラノ
審査請求日	平成19年12月7日(2007.12.7)	(72) 発明者	中臺 一博
(31) 優先権主張番号	60/633, 351		埼玉県和光市本町8-1、株式会社ホンダ
(32) 優先日	平成16年12月3日(2004.12.3)		・リサーチ・インスティテュート・ジャパン内
(33) 優先権主張国	米国 (US)	(72) 発明者	辻野 広司
早期審査対象出願			埼玉県和光市本町8-1、株式会社ホンダ
			・リサーチ・インスティテュート・ジャパン内

最終頁に続く

(54) 【発明の名称】 音声認識装置

(57) 【特許請求の範囲】

【請求項1】

外部から集音された音響信号から音声を認識するための音声認識装置であって、
前記音響信号を検出する少なくとも2つの音検出手段と、
前記音響信号に基づいて、音源の方向を求める音源定位手段と、
前記求められた音源の方向に基づいて音声分離する第1の手段と、
前記音声分離する第1の手段によって、前記分離の結果の信頼性に応じてマスクを生成する手段と、

前記音響信号の特徴量を抽出する手段と、

前記マスクを前記抽出された特徴量に適用して前記音響信号から音声認識する手段と

10

を備え、

前記マスクを生成する手段は、

音声分離する第1の手段で用いられる音源分離手段とは異なる音源分離法を用いて、
前記求められた音源の方向に基づいて、音響信号から音源に応じた音声分離する第2の手段と、

前記音声分離する第1の手段と前記音声分離する第2の手段によってなされた分離の結果を比較する手段と、

比較の結果に応じて音声のサブバンドにマスクした値を割り当てる手段と、

を備える、音声認識装置。

20

【請求項 2】

前記第 1 の音声を分離する手段は、
 音声の周波数サブバンドを定める手段を備え、
 前記サブバンドの位相差および音圧差の一方または両方が通過帯域内である、
 請求項 1 に記載の音声認識装置。

【請求項 3】

少なくとも 2 つの音検出によって集音された、音響信号を認識する方法であって、
 前記音響信号に基づいて音源を定位し、前記音源の方向を求めるステップと、
 前記求められた音源の方向に基づいて、音声を分離する第 1 のステップと、
 前記音声を分離する方法によって、分離の結果の信頼性に応じてマスクを生成するステ
 ップと、
 前記音響信号の特徴量を抽出するステップと、
 前記マスクを前記抽出された特徴量に適用して、前記音響信号から音声を認識するステ
 ップと、
 を含み、
 前記マスクを生成するステップは、
音声を分離する第 1 のステップで用いられる音源分離手段とは異なる音源分離法を用い
て、前記求められた音源の方向に基づいて、音響信号から音源に応じた音声を分離する第
2 のステップと、
前記音声を分離する第 1 のステップと前記音声を分離する第 2 のステップによってなさ
れた分離の結果を比較するステップと、
比較の結果に応じて音声のサブバンドにマスクした値を割り当てるステップと、
を含む、音響信号を認識する方法。

【請求項 4】

前記音声を分離する第 1 のステップは、
 音声の周波数サブバンドを定めるステップを含み、
 前記サブバンドの位相差および音圧差の一方または両方が通過帯域内である、
 請求項 3 に記載の音響信号を認識する方法。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、音声認識装置に関する。特に、雑音などによって劣化した音声に対し頑健な音声認識装置に関する。

【背景技術】

【0002】

一般に、実環境で利用される音声認識装置には、雑音や残響音の混入、入力装置の仕様などによって劣化した音声が入力される。この問題に対し、スペクトルサブトラクションやブラインド信号分離などの手法を用いて、音声認識の頑健さを向上させる取り組みが行われてきた。

【0003】

これらの取り組みの一環として、Sheffield大のM. Cookeらは、ミッシングフィーチャー理論を提案している(Martin Cooke, et al., "Robust automatic speech recognition with missing and unreliable acoustic data", SPEECH COMMUNICATION 34, p. 267-285, 2001を参照)。この手法は、入力音声の特徴量のうち、ミッシングフィーチャー(劣化した特徴量)を同定しマスクしてから認識することによって音声認識の頑健性向上を図るものであり、他の手法に比べて必要な雑音に関する知識が少ない、という特性を持つ。

【0004】

ミッシングフィーチャー理論において、劣化した特徴量の同定は、劣化していない音声の特徴量との差や、スペクトログラムの局所的なSN比、あるいはASA (Auditory Scene Analysis、聴覚情景分析)によって行われる。ASAは、スペクトルの調波構造やオンセットの

10

20

30

40

50

同期、音源の位置など、同じ音源から放射された音が共有する手掛かりを利用して、特徴量の要素をグループ化する方法である。音声認識は、マスクされた部分の元の特徴量を推定して認識する方法や、マスクされた特徴量に対応した音響モデルを生成して認識する方法などがある。

【発明の開示】

【発明が解決しようとする課題】

【0005】

ミッシングフィーチャー理論で音声認識の頑健性の向上を試みる場合、劣化した特徴量の同定に困難を伴うことが多い。本発明は、劣化した特徴量を完全に同定できない音声入力に対して音声認識の頑健性を向上させる音声認識装置を提案する。

10

【課題を解決するための手段】

【0006】

本発明は、外部から集音された音響信号から音声を認識するための音声認識装置を提供する。この装置は、音響信号を検出する少なくとも2つの音検出手段と、音響信号に基づいて音源の方向を求める音源定位部と、音源の方向に基づいて音響信号から音源による音声を分離する音源分離部と、分離の結果の信頼性に応じてマスクの値を生成するマスク生成部と、音響信号の特徴量を抽出する特徴抽出部と、マスクを特徴量に適用して音響信号から音声を認識する音声認識部と、を有する。

【0007】

本発明では、音源による音声を音響信号から分離した結果の信頼性に応じてマスクの値を生成するので、音声認識の頑健性を向上させることができる。

20

【0008】

本発明の一実施形態によると、マスク生成部が、音源分離部とは異なる複数の音源分離手法を用いて音響信号を分離した結果と、音源分離部による分離の結果との一致度合いに応じてマスクの値を生成する。

【0009】

本発明の一実施形態によると、マスク生成部が、音源方向によって定められる同一の音源かを判断するための通過幅に応じてマスクの値を生成する。

【0010】

本発明の一実施形態によると、複数の音源がある場合には、マスク生成部が該複数の音源のいずれか1つにだけ近いほど音源分離結果の信頼性を高めてマスクの値を生成する。

30

【発明を実施するための最良の形態】

【0011】

1. 概略

次に図面を参照して、この発明の実施の形態を説明する。図1は、本発明の一実施形態による音声認識装置10を含む音声認識システムを示す概略図である。

【0012】

図1に示すように、このシステムは、音声認識装置10を備えた躯体12が、その周囲にある音源14の発する音声を認識するものである。音源14は、人間やロボットなどコミュニケーション手段として音声を発するものである。躯体12は、移動ロボットや電

40

化製品など、インタフェースに音声認識を用いるものである。

【0013】

躯体12の両側には、音源からの音声を集音するための一对のマイク16a、16bが設置されている。なお、マイク16a、16bの位置は、躯体12の両側に限定されることなく、躯体12の他の位置に設置されても良い。また、マイクは、一对に限定されることなく、一对以上の個数が設置されても良い。

【0014】

このシステムは、音源14が発した音声を、マイク16を介して躯体12が集音する。集音された音声は躯体12内の音声認識装置10で処理される。音声認識装置10は、音声が発せられた音源14の方向を推定し、音声の内容を認識する。躯体12は例えば音声

50

の内容に応じたタスクを実施したり、自身の発話機構によって回答したりする。

【0015】

つづいて、音声認識装置10の詳細について説明する。図2は、本実施形態による音声認識装置10のブロック図である。

【0016】

複数のマイク16a、16bは、単数または複数の音源14が発した音声を集音し、これらの音声を含む音響信号を音声認識装置10に送る。

【0017】

音源定位部21は、マイク16a、16bより入力された音響信号から音源14の方向 θ を定位する。また、音源14や装置10自体が移動している場合は、定位された音源14の位置を時間方向に追跡する。本実施形態では、エピソード幾何、散乱理論、または伝達関数を利用して音源定位を実施する。

10

【0018】

音源分離部23は、音源定位部21で求められた音源14の方向情報 θ を利用し、入力信号から音源信号を分離する。本実施形態では、前述のエピソード幾何、散乱理論、または伝達関数を利用して得られるマイク間位相差 $\Delta\tau$ またはマイク間音圧差 Δp と、人間の聴覚特性を模した通過幅関数と、を組み合わせる音源分離を実施する。

【0019】

マスク生成部25は、音源分離部23の分離結果が信頼できるかどうかに応じて、マスクの値を生成する。信頼できるかどうかを求めるのに、入力信号のスペクトルや音源分離の結果を利用する。マスクは0~1の値をとり、1に近いほど信頼できる。マスク生成部で生成されたマスクの値はそれぞれ、音声認識に用いられる入力信号の特徴量に適用される。

20

【0020】

特徴抽出部27は、入力信号のスペクトルより特徴量を抽出する。

【0021】

音声認識部29は、音響モデルより特徴量の出力確率を求め、音声認識を行う。その際、マスク生成部25で生成したマスクを適用して、出力確率を調整する。本実施形態では、隠れマルコフモデル(Hidden Markov Model:HMM)によって認識を行う。

【0022】

以下、音声認識装置10の各構成要素で行われる処理について説明する。

30

【0023】

2. 音源定位部

音源定位部21は、複数のマイク16より入力された音響信号から音源14の方向を定位する。また、音源14や装置10自体が移動している場合は、定位された音源14の位置を時間方向に追跡する。本実施形態では、音源14およびマイク16のエピソード幾何を利用した音源定位(2.1節)、散乱理論を利用した音源定位(2.2節)、および伝達関数を利用した音源定位(2.3節)のうち1つを適用する。なお、音源定位の処理は、ビームフォーミングなど、その他の公知の方法を用いてもよい。

【0024】

2.1 音源およびマイクのエピソード幾何を利用した音源定位

この方法は、図3に示されるような、マイク16と音源14のエピソード幾何を利用して音源方向 θ を算出する。図3では、マイク16aおよびマイク16b間の距離は2bであり、両マイク間の中点を原点とし、原点から垂直方向を正面としている。

40

【0025】

なお、エピソード幾何の詳細については、例えば中臺他、“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”、AI Challenge研究会、pp. 1043-1049、人工知能学会、2002を参照されたい。

【0026】

エピソード幾何を利用した音源定位は、以下の手順で実施される。

50

【 0 0 2 7 】

1) マイク 1 6 a、1 6 b から入力された音響信号をFFTなどで周波数分析し、スペクトル $S1(f)$ 、 $S2(f)$ を求める。

2) 得られたスペクトルを複数の周波数領域(サブバンド)に分割し、各サブバンド f_i の位相差 $\Delta\phi(f_i)$ を、式(1)より求める。

【数1】

$$\Delta\phi(f_i) = \arctan\left(\frac{\text{Im}[S1(f_i)]}{\text{Re}[S1(f_i)]}\right) - \arctan\left(\frac{\text{Im}[S2(f_i)]}{\text{Re}[S2(f_i)]}\right) \quad (1)$$

10

ここで、 $\Delta\phi(f_i)$ は f_i のマイク間位相差である。 $\text{Im}[S1(f_i)]$ は、マイク1のサブバンド f_i におけるスペクトル $S1(f_i)$ の虚部であり、 $\text{Re}[S1(f_i)]$ は、マイク1のサブバンド f_i におけるスペクトル $S1(f_i)$ の実部である。 $\text{Im}[S2(f_i)]$ は、マイク2のサブバンド f_i におけるスペクトル $S2(f_i)$ の虚部であり、 $\text{Re}[S2(f_i)]$ は、マイク2のサブバンド f_i におけるスペクトル $S2(f_i)$ の実部である。

3) エピポーラ幾何(図3)を利用して式(2)を導出する。

【数2】

$$\Delta\phi(\theta, f_i) = \frac{2\pi f_i}{v} \times b(\theta + \sin \theta) \quad (2)$$

20

ここで、 v は音速を表し、 b は原点とマイクとの距離を表し、 θ は音源方向の角度を表す。

式(2)の θ に -90 度から $+90$ 度の範囲で例えば5度おきに代入して、図4に示すような周波数 f_i と位相差 $\Delta\phi$ との関係を求める。図4に示す関係を用いて、 $\Delta\phi(f_i)$ にもっとも近い θ (θ_i)を求め、この θ_i をサブバンド f_i の音源方向 θ_i とする。

4) 各サブバンドの音源方向 θ_i と周波数から、音源方向が近くかつ調音関係にあるものを選んでグループ化し、そのグループの音源方向 θ_s とする。なお、複数のグループが選別された場合、複数の音源が存在すると考えられるので、それぞれの音源方向を求めても良い。あらかじめ音源の数が分かっている場合は、音源の数に対応したグループ数を選ぶのが望ましい。

30

【 0 0 2 8 】

2.2 散乱理論を利用した音源定位

この方法は、マイク1 6を設置する躯体1 2による散乱波を考慮して、音源方向 θ_s を算出する。ここではマイク1 6を設置する躯体1 2をロボットの頭部とし、半径 b の球と仮定する。また、頭部の中心を極座標 (r, θ, ϕ) の原点とする。

【 0 0 2 9 】

なお、散乱理論の詳細については、例えばLax et al., "Scattering Theory", Academic Press, NY., 1989を参照されたい。

【 0 0 3 0 】

散乱理論を利用した音源定位は、以下の手順で実施される。

40

【 0 0 3 1 】

1) マイク1 6 a、1 6 b から入力された音響信号を、FFTなどで周波数分析しスペクトル $S1(f)$ 、 $S2(f)$ を求める。

2) 得られたスペクトルを複数の周波数領域(サブバンド)に分割し、各サブバンド f_i の位相差 $\Delta\phi(f_i)$ を、式(1)より求める。または、各サブバンド f_i の音圧差 $\Delta\rho(f_i)$ を、式(3)より求める。

【数3】

$$\Delta\rho(f_i) = 20 \log_{10} \frac{|P1(f_i)|}{|P2(f_i)|} \quad (3)$$

50

ここで、 (f_i) は両マイク間音圧差である。 $P1(f_i)$ はマイク 1 のサブバンド f_i のパワーであり、 $P2(f_i)$ はマイク 2 のサブバンド f_i のパワーである。

3) 音源 1 4 の位置を $r_0 = (r_0, 0, 0)$ 、観測点(マイク 1 6)の位置を $r = (b, 0, 0)$ 、音源と観測点の距離を $R = |r_0 - r|$ とすると、ロボット頭部における直接音によるポテンシャル V^i は、式(4)で定義される。

【数 4】

$$V^i = \frac{v}{2\pi R f} e^{i \frac{2\pi R f}{v}} \quad (4)$$

ここで、 f は周波数であり、 v は音速であり、 R は音源と観測点の距離である。

10

4) ロボット頭部における音源方向 (θ, f) からの直接音と散乱音によるポテンシャル $S(\theta, f)$ は、式(5)で定義される。

【数 5】

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

$$= - \left(\frac{v}{2\pi b f} \right)^2 \sum_{n=0}^{\infty} (2n+1) P_n(\cos \theta) \frac{h_n^{(1)} \left(\frac{2\pi r_0}{v} f \right)}{h_n^{(1)} \left(\frac{2\pi b}{v} f \right)} \quad (5)$$

20

ここで、 V^s は散乱音によるポテンシャルを表し、 P_n は第一種ルシャンドル (Legendre) 関数を表し、 $h_n(l)$ は第一種球ハンケル関数を表す。

5) マイク 1 6 a の極座標を $(b, \theta/2, 0)$ 、マイク 1 6 b の極座標を $(b, -\theta/2, 0)$ とすると、各マイクでのポテンシャルは、式(6)、(7)で表される。

$$S1(\theta, f) = S(\theta/2, f) \quad (6)$$

$$S2(\theta, f) = S(-\theta/2, f) \quad (7)$$

6) 音源の方向 (θ, f_i) と、各サブバンド f_i における位相差 (θ, f_i) 、音圧差 (θ, f_i) は、それぞれ式(8)、(9)によって関係付けられる。

【数 6】

$$\Delta\phi(\theta, f_i) = \arg(S1(\theta, f_i)) - \arg(S2(\theta, f_i)) \quad (8)$$

$$\Delta\rho(\theta, f_i) = 20 \log_{10} \frac{|S1(\theta, f_i)|}{|S2(\theta, f_i)|} \quad (9)$$

30

7) 予め式(8)、(9)の (θ, f_i) に適当な値(例えば5度毎)を入れ、周波数 f_i と位相差 (θ, f_i) との関係、または周波数 f_i と音圧差 (θ, f_i) との関係を求める。

8) (θ, f_i) または (θ, f_i) の中で、 (θ, f_i) または (θ, f_i) にもっとも近いものを、各サブバンド f_i の音源方向 θ_i とする。

9) 各サブバンドの音源方向 θ_i と周波数から、音源方向が近くかつ調音関係にあるものを選んでグループ化し、そのグループの音源方向 θ_s とする。なお、複数のグループが選別された場合、複数の音源が存在すると考えられるので、それぞれの音源方向を求めても良い。あらかじめ音源の数が分かっている場合は、音源の数に対応したグループ数を選ぶのが望ましい。また、 (θ, f_i) 、 (θ, f_i) の両方を使って音源方向 θ_s を求めてもよい。

40

【0032】

2.3 伝達関数を利用した音源定位

位相差や音圧差と周波数、音源方向を対応づけるのに一般的な方法は、伝達関数の測定である。伝達関数は、躯体 1 2 (たとえばロボット) に設置したマイク 1 6 a、1 6 b で、さまざまな方向からのインパルス応答を測定して作成される。これを用いて音源方向を

50

定位する。伝達関数を利用した音源定位は、以下の手順で実施される。

【 0 0 3 3 】

1) マイク 1 6 a、1 6 b から入力された音響信号を、FFTなどで周波数分析しスペクトル $S1(f)$ 、 $S2(f)$ を求める。

2) 得られたスペクトルを複数の周波数領域 (サブバンド) に分割し、各サブバンド f_i の位相差 (θ, f_i) を、式 (1) より求める。または、各サブバンド f_i の音圧差 (ρ, f_i) を、式 (3) より求める。

3) 適当な間隔 (例えば 5 度間隔) で ± 90 度の範囲で、インパルス応答を計測して伝達関数を取得する。方向 θ ごとにインパルス応答をマイク 1 6 a、1 6 b で測定してFFTなどの手法で周波数分析し、インパルス応答に対する各周波数 f のスペクトル (伝達関数) $Sp1(f)$ 、 $Sp2(f)$ を求める。伝達関数 $Sp1(f)$ 、 $Sp2(f)$ より、位相差 (θ, f) および音圧差 (ρ, f) を以下の式 (10)、式 (11) を用いて求める。

$$\Delta\phi(\theta, f) = \arg(Sp1(f)) - \arg(Sp2(f)) \quad (10)$$

$$\Delta\rho(\theta, f) = 20 \log_{10} \frac{|Sp1(f)|}{|Sp2(f)|} \quad (11)$$

± 90 度の範囲の任意の間隔の方向 θ と任意の周波数 f について上記計算を行い、算出された位相差 (θ, f) および音圧差 (ρ, f) の一例を図 5 および図 6 に示す。

4) 図 5 または図 6 に示す関係を用いて、 (θ, f_i) または (ρ, f_i) にもっとも近い θ_i を求め、それを各サブバンド f_i の音源方向 θ_i とする。

5) 各サブバンドの音源方向 θ_i と周波数から、音源方向が近くかつ調音関係にあるものを選んでグループ化し、そのグループの音源方向 θ_s とする。なお、複数のグループが選別された場合、複数の音源が存在すると考えられるので、それぞれの音源方向を求めても良い。また、 (θ, f_i) 、 (ρ, f_i) の両方を使って音源方向 θ_s を求めてもよい。

【 0 0 3 4 】

2.4 各マイクの入力信号の相互相関を利用した音源定位

この方法は、マイク 1 6 a、1 6 b の入力信号の相互相関から、音源 1 4 からマイク 1 6 a およびマイク 1 6 b への距離の差 (図 7 の d) を求め、マイク間距離 $2b$ との関係から音源方向 θ_s を推定する。この方法は以下の手順で実施される。

【 0 0 3 5 】

1) マイク 1 6 a およびマイク 1 6 b に入力された信号の相互相関 $CC(T)$ を式 (11) で計算する。

【数 8】

$$CC(T) = \int_0^T x_1(t)x_2(t+T)dt \quad (12)$$

ここで、 T はフレーム長を表す。 $x_1(t)$ はフレーム長 T で切り出されたマイク 1 6 a からの入力信号を表し、 $x_2(t)$ はフレーム長 T で切り出されたマイク 1 6 b からの入力信号を表す。

2) 得られた相互相関からピークを抽出する。抽出するピーク数は、あらかじめ音源数が分かっている場合は、音減数と同数を抽出するのが望ましい。抽出したピークの時間軸上の位置が、マイク 1 6 a およびマイク 1 6 b への信号の到達時間差を示す。

3) 信号の到達時間差と音速より、音源 1 4 からマイク 1 6 a、1 6 b までの距離の違い (図 7 の d) を算出する。

4) 図 7 に示すように、マイク間距離 $2b$ および音源からマイクへの距離の差 d を用いて、式 (12) から音源 1 4 の方向 θ_s を求める。

$$\theta_s = \arcsin(d/2b) \quad (13)$$

複数のピークを抽出した場合は、それぞれ音源方向 θ_s を求める。

10

20

30

40

50

【 0 0 3 6 】

2.5 音源方向追跡

音源 1 4 または 躯体 1 2 が移動する場合には、音源方向の追跡を行う。図 8 は、音源方向 θ_s の時間変化を示す。追跡は、それまでの時刻で得られた θ_s の軌跡から予測される音源方向 θ_p と、実際に得られた θ_s とを比較し、その差が予め定めたしきい値よりも小さい場合には、同一音源からの信号と判断し、しきい値よりも大きい場合は、同一音源からの信号ではないと判断して行う。予測には、カルマンフィルタや自己回帰予測、HMM等、既存の時系列信号予測手法を用いる。

【 0 0 3 7 】

3. 音源分離部

音源分離部 2 3 は、音源定位部 2 1 で求められた音源 1 4 の方向情報 θ_s を利用し、入力信号から音源信号を分離する。本実施形態では、前述のエピポーラ幾何、散乱理論、または伝達関数を利用して得られるマイク間位相差 $\Delta\theta$ またはマイク間音圧差 Δp と、人間の聴覚特性を模した通過幅関数 $H(\theta)$ と、を組み合わせた分離方法について述べる。しかし、音源分離部 2 3 で用いる手法は、ビームフォーミングやGSS (Geometric Source Separation、幾何学的信号源分離) など、音源方向を利用し、かつサブバンドごとに音源分離をする公知の手法を用いてもよい。音源分離が時間領域で行われる場合は、分離の後周波数領域に変換する。本実施形態では音源分離は以下の手順で実施される。

【 0 0 3 8 】

1) 音源定位部 2 1 より音源方向 θ_s と、入力信号のスペクトルのサブバンド f_i の位相差 $\Delta\theta(f_i)$ または音圧差 $\Delta p(f_i)$ を受け取る。音源分離部 2 3 で周波数領域における音源定位の手法を用いない場合には、ここで式 (1) または式 (3) を用いて $\Delta\theta(f_i)$ または $\Delta p(f_i)$ を求める。

2) 音源方向と通過幅の関係を示す通過幅関数を用いて、音源定位部 2 1 で得られた音源方向 θ_s に対応する通過幅 $H(\theta_s)$ を求める。

通過幅関数は、音源方向に対する解像度が正面方向では高く周辺では低いという人の聴覚特性に基づき設計された関数であり、例えば図 9 に示すように正面方向の通過幅が狭く、周辺の通過幅が広がっている。横軸は、躯体 1 2 の正面を 0 [deg] とした場合の水平角である。

3) 得られた $H(\theta_s)$ より、通過帯域の下限 f_l と上限 f_h (図 8 に例示) を、式 (10) を用いて算出する。

【数 9】

$$\begin{aligned}\theta_l &= \theta_s - \delta(\theta_s) \\ \theta_h &= \theta_s + \delta(\theta_s)\end{aligned}\quad (14)$$

4) f_l 、 f_h に対応する位相差 $\Delta\theta(f_l)$ 、 $\Delta\theta(f_h)$ を、前述のエピポーラ幾何 (式 (2) および図 4)、散乱理論 (式 (8))、伝達関数 (図 5) のいずれかを用いて推定する。図 11 は推定した位相差と周波数 f_i との関係の一例を示すグラフである。または、 f_l 、 f_h に対応する音圧差 $\Delta p(f_l)$ 、 $\Delta p(f_h)$ を、前述の散乱理論 (式 (9))、伝達関数 (図 6) のいずれかを用いて推定する。図 12 は推定した音圧差と周波数 f_i との関係の一例を示すグラフである。

5) 各サブバンドの $\Delta\theta(f_i)$ または $\Delta p(f_i)$ が、通過帯域内にあるかどうか調べ、通過帯域内のものを選択する (図 11、図 12)。一般に、低周波数の定位は位相差、高周波数の定位は音圧差を利用するほうが、分離精度が増すと言われているので、予め定めたしきい値 (例えば 1500 [Hz]) より小さいサブバンドは位相差 $\Delta\theta$ を、大きいサブバンドは音圧差 Δp を使って選択しても良い。

6) 選択されたサブバンドのフラグを 1 に設定し、選択されなかったサブバンドのフラグを 0 に設定する。1 のフラグがついたサブバンドが、音源信号として分離される。

【 0 0 3 9 】

なお、音源分離を、今まで述べてきた線形周波数領域のスペクトルではなく、メル周波

10

20

30

40

50

数領域のスペクトルで行ってもよい。メル周波数とは、音の高低に対する人間の間隔尺度であり、その値は実際の周波数の対数にほぼ対応する。この場合は、前述の音源分離部 2 3 の処理のステップ 1) の後に、メル周波数に変換するフィルタ処理を加えた以下の手順で、メル周波数領域での音源分離を行う。

【 0 0 4 0 】

1) マイク 1 6 a、1 6 b に入力された信号を、FFTなどで周波数分析しスペクトル $S_1(f)$ 、 $S_2(f)$ を求める。

2) メル周波数領域で等間隔に配置した三角窓 (例えば 2 4 個) によりフィルタバンク分析を行う。

3) 得られたメル周波数領域スペクトルの各サブバンド m_j の位相差 (m_j) を式 (1) (ただし f_j 、 m_j) より求める。またはマイク間音圧差 (m_j) を、式 (3) (ただし f_j 、 m_j) より求める。 10

4) 音源方向と通過幅の関係を示す通過幅関数 (図 9) を用いて、音源定位部 2 1 で得られた音源方向 s に対応する通過幅 (s) を求める。

5) 得られた (s) より、通過帯域の下限 l と上限 h を、式 (1 0) を用いて算出する。

6) l 、 h に対応する位相差 l 、 h を、前述のエピポーラ幾何 (式 (2) および図 4)、散乱理論 (式 (8))、伝達関数 (図 5) のいずれかを用いて推定する。または、 l 、 h に対応する音圧差 l 、 h を、前述の散乱理論 (式 (9))、伝達関数 (図 6) のいずれかを用いて推定する。 20

7) 各メル周波数の (m_j) または (m_j) が、通過帯域内にあるかどうか調べ、通過帯域内のものを選択する。一般に、低周波数の定位は位相差、高周波数の定位は音圧差を利用するほうが、分離精度が増すと言われているので、予め定めたいきい値 (例えば 1500 [Hz]) より小さいサブバンドは位相差 を、大きいサブバンドは音圧差 を使って選択しても良い。

8) 選択されたメル周波数に 1 のフラグを設定し、選択されなかったメル周波数に 0 のフラグを設定する。1 のフラグがついたメル周波数を分離された信号とする。

【 0 0 4 1 】

なお、音源分離がメル周波数領域で求められた場合、後述するマスク生成部 2 5 で行われるメル周波数への変換は不要となる。 30

【 0 0 4 2 】

4 . マスク生成部

マスク生成部 2 5 は、音源分離部 2 3 の分離結果が信頼できるかどうかに応じて、マスクの値を生成する。本実施形態では、複数の音源分離方法からの情報を利用したマスク生成 (4 . 1 節)、通過幅関数を利用したマスク生成 (4 . 2 節)、複数音源の影響を考慮したマスク生成 (4 . 3 節) のいずれかを適用する。音源分離部 2 3 で設定されたフラグ (0 または 1) の信頼度を調べ、フラグの値と信頼度を考慮してマスクの値を設定する。マスクは 0 ~ 1 の値をとり、1 に近いほど信頼できるものとする。

【 0 0 4 3 】

4 . 1 複数の音源分離方法からの情報を利用したマスク生成 40

ここでは、複数の音源分離方法による信号分離の結果を用いて、音源分離部 2 3 の分離結果が信頼できるかどうかを確認し、マスクを生成する。この処理は以下の手順で実施される。

【 0 0 4 4 】

1) 音源分離部 2 3 で用いられていない音源分離手法を少なくとも 1 つ用いて音源分離を行い、音源分離部 2 3 と同様にサブバンドごとにフラグを立てる。本実施形態では、音源分離部 2 3 では以下の要素のいずれかを用いて音源分離が実施される。

- i) エピポーラ幾何に基づく位相差
- ii) 散乱理論に基づく位相差
- iii) 散乱理論に基づく音圧差 50

- iv) 伝達関数に基づく位相差
v) 伝達関数に基づく音圧差

【0045】

2) 音源分離部23で得られたフラグと、1)で得られたフラグのそれぞれが一致しているかどうかを調べ、マスクを生成する。例えば、音源分離部23の手法にi)エピソード幾何に基づく位相差を用い、マスク生成部25の手法にii)散乱理論に基づく位相差、iii)散乱理論に基づく音圧差、およびv)伝達関数に基づく音圧差を用いる場合を考えると、各状態におけるマスクの値は以下のようになる。

【表1】

i) のフラグ	ii) iii) v) のフラグ	マスク値
0	全て0	1
0	2個が0	1/3
0	1個以下が0	0
1	全て1	1
1	2個が1	1/3
1	1個以下が1	0

10

20

【0046】

3) 得られたマスク値を、メルスケールのフィルタバンク分析を行って、メル周波数軸に変換し、マスクを生成する。なお、上述のように、音源分離がメル周波数領域で求められた場合には、このステップは不要である。

【0047】

また、メル周波数軸に変換したマスクの値に対して適当なしきい値を設けておき、しきい値を超えたものは1、そうでないものは0をとる二値マスクに変換してもよい。

30

【0048】

4.2 通過幅関数を利用したマスク生成

この方法では、音源方向 θ_s と通過幅関数 (g_s) を利用し、音源方向との近さによってマスク値を生成する。つまり、音源方向に近いほど、音源分離部23で付された1のフラグは信頼でき、音源方向から遠いほど、音源分離部23で付された0のフラグは信頼できると考える。この処理は以下の手順で実施される。

【0049】

- 1) 音源定位部21より、音源方向 θ_s と入力信号を受け取る。

【0050】

2) 入力信号より、各サブバンドの音源方向 θ_i を求める(音源定位部21で音源方向が求められている場合は、それを利用する)。

40

【0051】

3) 音源分離部23より、通過幅 (g_s) と各サブバンド f_i のフラグを受け取る(以下 f_i とする)。

【0052】

4) f_i を用いてマスクの関数を生成し、各サブバンドの θ_i と比べて仮マスクを求める。関数は次式のように与えられ、図13に示すような挙動となる。

【数 10】

$$\text{仮マスク} = \begin{cases} 1 & (-\pi \leq \theta_i < \theta_s - 2\theta_t) \\ -\frac{\theta_i - \theta_s}{\theta_t} - 1 & (\theta_s - 2\theta_t \leq \theta_i < \theta_s - \theta_t) \\ \frac{\theta_i - \theta_s}{\theta_t} + 1 & (\theta_s - \theta_t \leq \theta_i < \theta_s) \\ -\frac{\theta_i - \theta_s}{\theta_t} + 1 & (\theta_s \leq \theta_i < \theta_s + \theta_t) \\ \frac{\theta_i - \theta_s}{\theta_t} - 1 & (\theta_s + \theta_t \leq \theta_i < \theta_s + 2\theta_t) \\ 1 & (\theta_s + 2\theta_t \leq \theta_i < \pi) \end{cases} \quad (15)$$

10

【0053】

5) 音源分離部 23 で求めたサブバンド f_i のフラグと、ステップ 4) で求めた仮マスクから、以下の通りマスクを生成する。

【表 2】

フラグ	仮マスク	マスク値
0	1	<u>1</u>
0	$1 > \text{仮マスク} > 0$	仮マスクの値
0	0	<u>0</u>
1	1	1
1	$1 > \text{仮マスク} > 0$	仮マスクの値
1	0	0

20

30

【0054】

6) 得られたマスク値を、メルスケールのフィルタバンク分析を行って、メル周波数軸に変換し、マスクを生成する。なお、上述のように、音源分離がメル周波数領域で求められた場合には、このステップは不要である。

【0055】

また、メル周波数軸に変換したマスクの値に対して適当なしきい値を設けておき、しきい値を超えたものは 1、そうでないものは 0 をとる二値マスクに変換してもよい。

【0056】

4.3 複数音源の影響を考慮したマスク生成

ここでは、音源が複数ある場合に、2つ以上の音源の信号が含まれていると推定されるサブバンドの信頼性を下げるように、マスク値を生成する。

【0057】

1) 音源定位部 21 より、音源方向 s_1, s_2, \dots と入力信号を受け取る。

【0058】

2) 入力信号より、各サブバンドの音源方向 θ_i を求める。音源定位部 21 で音源方向が求められている場合は、それを利用する。

【0059】

3) 音源分離部 23 より、各音源方向 s_1, s_2, \dots の通過帯域 (f_{l1}, f_{h1})、(f_{l2}, f_{h2})、...

50

, h_2), ... とフラグを受け取る。

【0060】

4) 各サブバンドの音源方向 θ_i が、

- i) 2つ以上の音源の通過帯域 (f_l, f_h) に含まれている
- ii) その音源の通過帯域にも含まれていない

かどうか調べ、i) または ii) にあてはまるサブバンドには 0、それ以外には 1 の仮マスクを生成する。

【0061】

5) フラグと仮マスクより、以下の通りマスクを生成する。

【表3】

フラグ	仮マスク	マスク値
0	1	0
0	0	1
1	1	1
1	0	0

10

20

【0062】

6) 得られたマスク値を、メルスケールのフィルタバンク分析を行って、メル周波数軸に変換し、マスクを生成する。なお、上述のように、音源分離がメル周波数領域で求められた場合には、このステップは不要である。

【0063】

また、メル周波数軸に変換したマスクの値に対して適当なしきい値を設けておき、しきい値を超えたものは 1、そうでないものは 0 をとる二値マスクに変換してもよい。

【0064】

5. 特徴抽出部

特徴抽出部 27 は、一般的に知られる手法を用いて、入力信号のスペクトルより特徴量を求める。この処理は以下の手順で実施される。

1) FFT等でスペクトルを求める。

2) メル周波数領域で等間隔に配置した三角窓 (例えば24個) によりフィルタバンク分析を行う。

3) 分析結果の対数を取り、メル周波数対数スペクトルを得る。

4) 対数スペクトルを離散コサイン変換する。

5) ケプストラム係数の0次と高次 (例えば13次から23次) の項を 0 にする。

6) ケプストラム平均除去を行う。

7) 逆離散コサイン変換を行う。

【0065】

以下、求められた特徴量を、特徴ベクトル $x = (x_1, x_2, \dots, x_j, \dots, x_J)$ として扱う。

【0066】

6. 音声認識部

本実施形態では、音声認識部 29 は、従来技術として知られるHMMによって音声認識を行う。

【0067】

特徴ベクトル x 、状態 S の時の通常の連続分布型HMMの出力確率 $f(x, S)$ は、式 (16) で表される。

30

40

50

【数 1 1】

$$f(\mathbf{x} | S) = \sum_{k=1}^N P(k | S) f(\mathbf{x} | k, S) \quad (16)$$

ここで、Nは混合正規分布の混合数を表し、P(k|S)は混合比を表す。

【0068】

ミッシングフィーチャー理論に基づく音声認識では、f(x, S)をxの確率密度関数p(x)で平均したものを利用する。

【数 1 2】

$$\overline{f(\mathbf{x} | S)} = \sum_{k=1}^N P(k | S) f(\mathbf{x}_r | k, S) \quad (17)$$

ここで、 $\mathbf{x} = (\mathbf{x}_r, \mathbf{x}_u)$ とし、 \mathbf{x}_r は特徴ベクトルのうち信頼できる成分で、マスクが0より大きいもの、 \mathbf{x}_u は特徴ベクトルのうち信頼できない成分で、マスクが0のものを示す。

【0069】

信頼できない特徴成分が[0, \mathbf{x}_u]の範囲に一様分布すると仮定すると、式(17)は、式(18)に書き直せる。

【数 1 3】

$$\overline{f(\mathbf{x} | S)} = \sum_{k=1}^N P(k | S) f(\mathbf{x}_r | k, S) \frac{1}{\mathbf{x}_u} \int_0^{\mathbf{x}_u} f(\mathbf{x}'_r | k, S) d\mathbf{x}'_u \quad (18)$$

【0070】

xのj番目の成分の出力確率 $o(x_j | S)$ は、式(19)のように表せる。

【数 1 4】

$$o(x_j | S) = \begin{cases} M(j) f(x_j | S) + (1 - M(j)) \overline{f(x_j | S)} & \text{if } M(j) \neq 0 \\ 1 & \text{otherwise} \end{cases} \quad (19)$$

ここで、M(j)は特徴ベクトルのj番目の成分のマスクを表す。

【0071】

全体の出力確率 $o(x | S)$ は、式(20)のように表せる。

【数 1 5】

$$o(x | S) = \prod_{j=1}^J o(x_j | S) \quad (20)$$

ここでJは特徴ベクトルの次元を表す。

【0072】

式(20)は、式(21)でも表せる。

【数 1 6】

$$o(x | S) = \sum_{k=1}^N P(k | S) \exp \left\{ \sum_{j=1}^J M(j) \log f(x_j | k, S) \right\} \quad (21)$$

式(20)または式(21)を用いて音声認識を行う。

【0073】

以上にこの発明を特定の実施例について説明したが、この発明はこのような実施例に限定されるものではない。

10

20

30

40

50

【図面の簡単な説明】

【0074】

【図1】本発明の一実施形態による音声認識装置を含む音声認識システムを示す概略図である。

【図2】本実施形態による音声認識装置のブロック図である。

【図3】マイクおよび音源のエピポーラ幾何を示す図である。

【図4】エピポーラ幾何から導かれたマイク間位相差、周波数 f および音源方向 θ_s の関係を示す図である。

【図5】伝達関数から導かれたマイク間位相差、周波数 f 、および音源方向 θ_s の関係を示す図である。

10

【図6】伝達関数から導かれたマイク間音圧差、周波数 f 、および音源方向 θ_s の関係を示す図である。

【図7】マイクおよび音源の位置関係を示す図である。

【図8】音源方向 θ_s の時間変化を示す図である。

【図9】通過幅関数()を示す図である。

【図10】音源方向 θ_s と通過帯域を示す図である。

【図11】音源分離部における位相差によるサブバンド選択を示す図である。

【図12】音源分離部における音圧差によるサブバンド選択を示す図である。

【図13】通過幅関数を利用したマスクの関数を示す図である。

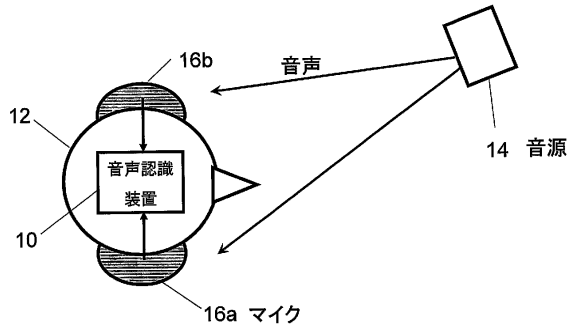
20

【符号の説明】

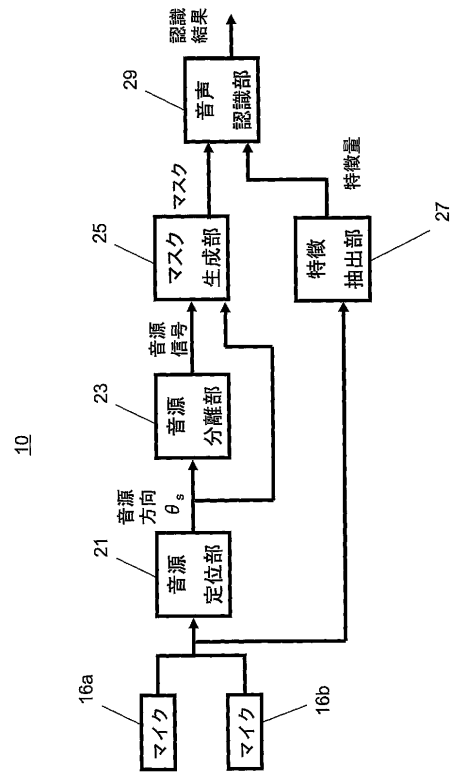
【0075】

- 10 音声認識装置
- 14 音源
- 16 マイク
- 21 音源定位部
- 23 音源分離部
- 25 マスク生成部
- 27 特徴抽出部
- 29 音声認識部

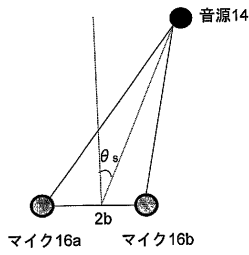
【図1】
第1図



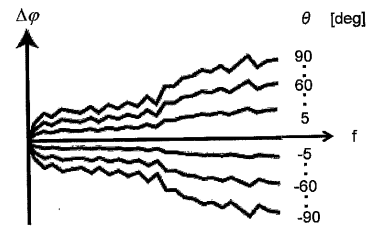
【図2】
第2図



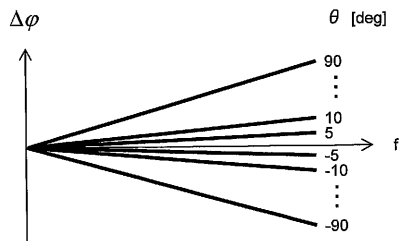
【図3】
第3図



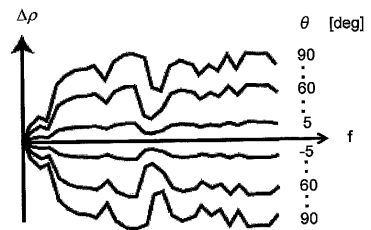
【図5】
第5図



【図4】
第4図

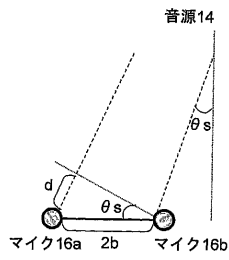


【図6】
第6図



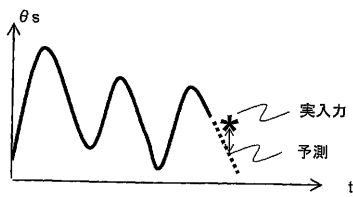
【図7】

第7図



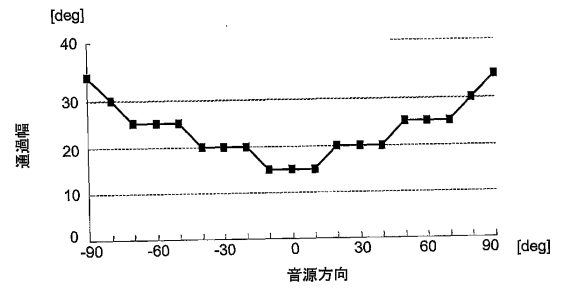
【図8】

第8図



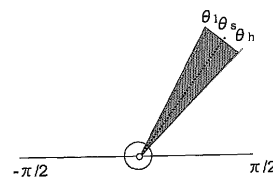
【図9】

第9図



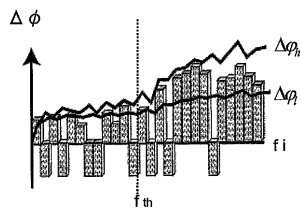
【図10】

第10図



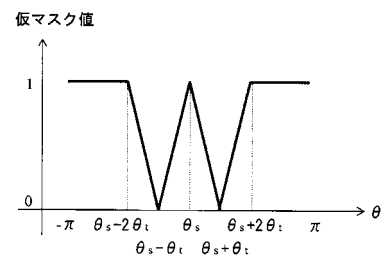
【図11】

第11図



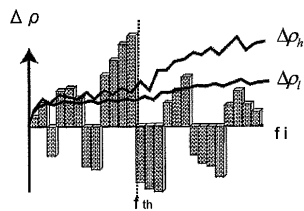
【図13】

第13図



【図12】

第12図



フロントページの続き

(72)発明者 奥乃 博

京都府京都市中京区東洞院通三条下る三文字町205番地の3-1102

(72)発明者 山本 俊一

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 Yamamoto, S.; Nakadai, K.; Tsujino, H.; Okuno, H.G., Assessment of general applicability of robot audition system by recognizing three simultaneous speeches, IEEE/RSJ International Conference on Intelligent Robots and Systems, 2004. (IROS 2004). Proceedings. 2004, 米国, IEEE, 2004年9月28日, vol.3, 2111-2116

Rickard, S.; Yilmaz, Z., On the approximate W-disjoint orthogonality of speech, IEEE International Conference on Acoustics, Speech, and Signal Processing, 2002. Proceedings. (ICASSP '02), 米国, IEEE, 2002年5月13日, vol.1, 2002, 1, 529-532

(58)調査した分野(Int.Cl., DB名)

G10L 15/00-15/28

G10L 21/00-21/06

H04R 3/00- 3/14