

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5180928号
(P5180928)

(45) 発行日 平成25年4月10日(2013.4.10)

(24) 登録日 平成25年1月18日(2013.1.18)

(51) Int.Cl. F I
G 1 O L 15/20 (2006.01) G 1 O L 15/20 3 7 O E
 G 1 O L 15/20 3 6 O B

請求項の数 9 (全 22 頁)

(21) 出願番号	特願2009-185164 (P2009-185164)	(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
(22) 出願日	平成21年8月7日(2009.8.7)	(74) 代理人	110000246 特許業務法人O F H特許事務所
(65) 公開番号	特開2010-49249 (P2010-49249A)	(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
(43) 公開日	平成22年3月4日(2010.3.4)	(72) 発明者	高橋 徹 京都府京都市左京区吉田本町 国立大学法 人京都大学 大学院情報学研究科内
審査請求日	平成23年11月24日(2011.11.24)	(72) 発明者	奥乃 博 京都府京都市左京区吉田本町 国立大学法 人京都大学 大学院情報学研究科内 最終頁に続く
(31) 優先権主張番号	61/136, 225		
(32) 優先日	平成20年8月20日(2008.8.20)		
(33) 優先権主張国	米国 (US)		

(54) 【発明の名称】 音声認識装置及び音声認識装置のマスク生成方法

(57) 【特許請求の範囲】

【請求項1】

複数音源からの混合音を分離する音源分離部と、
 前記音源分離部が分離を行った際の分離信頼度に対応して、分離された音声ごとに、0から1の間の連続的な値をとりうるソフトマスクを生成するマスク生成部と、
 前記音源分離部によって分離された音声を、前記マスク生成部で生成されたソフトマスクを使用して認識する音声認識部と、
 を備え、
 前記分離信頼度は、前記音源分離部により分離された音声毎に算出される、他の音源からの混ざり込みの程度を表わす数値であって、他の音源からの混ざり込みがなく完全に分離できている場合には1となり、混ざり込みが大きくなるにつれて0に近い値をとり、
 前記マスク生成部は、前記算出された分離信頼度のヒストグラムに基づいて前記ソフトマスクを生成する、
 音声認識装置。

【請求項2】

前記ソフトマスクは、前記ヒストグラムから算出される、前記分離信頼度の確率分布を構成する2つの正規分布の確率密度関数に基づいて定められる、

請求項1に記載の音声認識装置

【請求項3】

前記ソフトマスクが、Rを分離信頼度、a、bを定数として、Rのシグモイド関数

$1 / (1 + \exp(-a(R - b)))$
を使用して定められ、

上記定数 a 及び b は、前記 2 つの正規分布の確率密度関数に基づいて定められる、請求項 2 に記載の音声認識装置。

【請求項 4】

音声認識装置のソフトマスクを生成する方法であって、前記音声認識装置は、
複数音源からの混合音を分離する音源分離部と、
前記音源分離部が分離を行った際の分離信頼度に対応して、分離された音声ごとに、0
から 1 の間の連続的な値をとりうるソフトマスクを生成するマスク生成部と、

前記音源分離部によって分離された音声を、前記マスク生成部で生成されたソフトマスク
を使用して認識する音声認識部と、を備え、前記ソフトマスクは前記分離信頼度の関数
を使用して定められており、

分離信頼度のヒストグラムを求めるステップと、

分離信頼度のヒストグラムに基づいて、前記関数が有する少なくとも一つのパラメータ
の値を定めるステップと、を含み、

前記分離信頼度は、前記音源分離部により分離された音声毎に算出される、他の音源から
の混ざり込みの程度を表わす数値であって、他の音源からの混ざり込みがなく完全に分
離できている場合には 1 となり、混ざりこみが大きくなるにつれて 0 に近い値をとる、

音声認識装置のソフトマスクを生成する方法。

【請求項 5】

前記関数が有する他の少なくとも一つのパラメータの探索範囲を定めるステップと、
前記定められた探索範囲内において、前記他の少なくとも一つのパラメータの値を変化
させながら、前記音声認識装置の音声認識率を求めるステップと、

前記音声認識率が最大となる値を前記他の少なくとも一つのパラメータの値とするステ
ップとを含む、

請求項 4 に記載の音声認識装置のソフトマスクを生成する方法。

【請求項 6】

μ_1 、 μ_2 ($\mu_1 < \mu_2$) を平均値、 σ_1 、 σ_2 を標準偏差とし、分離信頼度を R とし
て、分離信頼度 R のヒストグラムを、 (μ_1, σ_1) を有する第 1 の正規分布の確率密度
関数 $f_1(R)$ 及び (μ_2, σ_2) を有する第 2 の正規分布の確率密度関数 $f_2(R)$ で
フィッティングすることによって、 μ_1 、 μ_2 、 σ_1 及び σ_2 を推定し、 $f_1(R)$ 、 $f_2(R)$ 、 μ_1 及び μ_2 を使用して前記ソフトマスクを生成する、請求項 4 に記載の音声
認識装置のソフトマスクを生成する方法。

【請求項 7】

前記ソフトマスクの値を $S(R)$ 、 $f(R) = f_1(R) + f_2(R)$ として、

$R < \mu_1$ において $S(R) = 0$

$\mu_1 < R < \mu_2$ において $S(R) = f_2(R) / f(R)$

$R > \mu_2$ において $S(R) = 1$

とする、請求項 6 に記載の音声認識装置のソフトマスクを生成する方法。

【請求項 8】

前記ソフトマスクの値を $S(R)$ 、

$R < \mu_1$ において

【数 1】

$$f_1'(R) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

$\mu_1 < R$ において

10

20

30

40

【数 2】

$$f1'(R) = f1(R)$$

R < μ 2 において

【数 3】

$$f2'(R) = f2(R)$$

μ 2 R において

【数 4】

$$f2'(R) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

とし、

【数 5】

$$f'(R) = f1'(R) + f2'(R)$$

として、

【数 6】

$$SM(R) = \frac{f2'(R)}{f'(R)}$$

とする、請求項 6 に記載の音声認識装置のソフトマスクを生成する方法。

【請求項 9】

f 1 (R) と f 2 (R) との交点で

$$\mu 1 < R < \mu 2$$

を満たす R の値を b とし、

$$1 / (1 + \exp (- a (R - b)))$$

が

$$f 2 (R) / f (R)$$

とフィッティングするように a を定めて、前記ソフトマスクの値を S (R) として、

$$S (R) = 1 / (1 + \exp (- a (R - b)))$$

とする、請求項 6 に記載の音声認識装置のソフトマスクを生成する方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数音源の音声を同時認識する音声認識装置及び音声認識装置のマスク生成方法に関する。

【背景技術】

【0002】

複数音源の音声を同時認識する技術は、たとえば、ロボットが実環境で活動する際に重要な技術である。複数音源の音声を同時認識する音声認識システムは、音源ごとに音声を分離し、分離した音声の音響特徴量を使用して音声認識を行なう。ここで、音声認識を行なう際に、分離の信頼度に応じて音響特徴量ごとにマスクが使用される（たとえば、非特許文献 1）。このようなマスクとしては、従来、0 または 1 の 2 値のハードマスクが使用されていた（たとえば、非特許文献 2）。0 から 1 の連続的な値を与えるソフトマスクも知られてはいたが（たとえば、非特許文献 3）、複数音源の音声を同時認識する音声認識

10

20

30

40

50

システム用のソフトマスクは開発されていなかった。その理由は、従来、当業者は、複数音源の音声を同時認識する音声認識にはハードマスクの方が適していると考えていたためである（たとえば、非特許文献2）。このように、複数音源の音声を同時認識する音声認識に適したソフトマスクを備え、音声認識率を向上させた音声認識装置は開発されていなかった。

【先行技術文献】

【非特許文献】

【0003】

【非特許文献1】M. L. Seltzer, B. Raj, and R. M. Stern, "A Bayesian frame work for spectrographic mask estimation for missing feature speech recognition," Spe 10
ech Communication, vol.43, pp. 379-393, 2004

【非特許文献2】Shun'ichi Yamamoto, Jean-Marc Valin, Kazuhiro Nakadai, Jean Rou at, Francois Michaud, Tetsuya Ogata, and Hiroshi G. Okuno, "Enhanced Robot Speec h Recognition Based on Microphone Array Source Separation and Missing Feature Th eory," in Proc. of IEEE CRA-2005, pp. 1489-1494, 2005

【非特許文献3】J. Barker, L. Josifovski, M. P. Cooke and P. D. Green, "Soft de cision in missing data techniques for robust automatic speech recognition," Pro c., ICSLP-2000, 2000

【発明の概要】

【発明が解決しようとする課題】

20

【0004】

したがって、複数音源の音声を同時認識する音声認識に適したソフトマスクを備え、音 声認識率を向上させた音声認識装置に対するニーズがある。

【課題を解決するための手段】

【0005】

本発明の音声認識装置は、複数音源からの混合音を分離する音源分離部と、前記音源分 離部が分離を行った際の分離信頼度に対応して、分離された音声ごとに、0から1の間の 連続的な値をとりうるソフトマスクを生成するマスク生成部と、前記音源分離部によって 分離された音声を、前記マスク生成部で生成されたソフトマスクを使用して認識する音声 認識部と、を備えている。

30

【0006】

本発明による音声認識装置によれば、分離信頼度に対応して、分離された音声ごとに、 生成された0から1の間の連続的な値をとりうるソフトマスクを使用して音声認識される ので、音声認識率が向上する。

【0007】

本発明の実施形態による音声認識装置においては、前記ソフトマスクが、Rを分離信頼 度、a、bを定数として、Rのシグモイド関数

$$1 / (1 + \exp (- a (R - b)))$$

を使用して定められている。

【0008】

40

本実施形態によれば、シグモイド関数の定数a及びbを変化させることにより、容易に ソフトマスクの調整を行うことができる。

【0009】

本発明の実施形態による音声認識装置においては、前記ソフトマスクが、Rを分離信頼 度として、Rを変数とする正規分布の確率密度関数を使用して定められている。

【0010】

本実施形態によれば、正規分布の確率密度関数の形状を変化させることにより、容易に ソフトマスクの調整を行うことができる。

【0011】

本発明による音声認識装置のソフトマスクを生成する方法は、複数音源からの混合音を

50

分離する音源分離部と、前記音源分離部が分離を行った際の分離信頼度に対応して、分離された音声ごとに、0から1の間の連続的な値をとりうるソフトマスクを生成するマスク生成部と、前記音源分離部によって分離された音声を、前記マスク生成部で生成されたソフトマスクを使用して認識する音声認識部と、を備えた音声認識装置のソフトマスクを生成する。前記ソフトマスクは、少なくとも一つのパラメータを有する分離信頼度の関数を使用して定められている。該方法は、前記少なくとも一つのパラメータの探索範囲を定めるステップと、前記少なくとも一つのパラメータの探索範囲内において、前記少なくとも一つのパラメータの値を変化させながら、前記音声認識装置の音声認識率を求めるステップと、前記音声認識率が最大となる値を前記少なくとも一つのパラメータの値とするステップとを含む。

10

【0012】

本発明による音声認識装置のソフトマスクを生成する方法によれば、前記ソフトマスクは、少なくとも一つのパラメータを有する分離信頼度の関数を使用して定められているので、少なくとも一つのパラメータの値を変化させながら、音声認識装置の音声認識率を求めことにより、確実に、音声認識率が最大となるように少なくとも一つのパラメータの値を定めることができる。

【0013】

本発明による音声認識装置のソフトマスクを生成する方法は、複数音源からの混合音を分離する音源分離部と、前記音源分離部が分離を行った際の分離信頼度に対応して、分離された音声ごとに、0から1の間の連続的な値をとりうるソフトマスクを生成するマスク生成部と、前記音源分離部によって分離された音声を、前記マスク生成部で生成されたソフトマスクを使用して認識する音声認識部と、を備えた音声認識装置のソフトマスクを生成する。前記ソフトマスクは、少なくとも一つのパラメータを有する分離信頼度の関数を使用して定められている。該方法は、分離信頼度のヒストグラムを求めるステップと、分離信頼度のヒストグラムの形状から前記少なくとも一つのパラメータの値を定めるステップと、を含む。

20

【0014】

本発明による音声認識装置のソフトマスクを生成する方法によれば、前記ソフトマスクは、少なくとも一つのパラメータを有する分離信頼度の関数を使用して定められているので、分離信頼度のヒストグラムを求めることにより、分離信頼度のヒストグラムの形状から適切少なくとも一つのパラメータの値を定めることができる。

30

【0015】

本発明の実施形態による音声認識装置のソフトマスクを生成する方法においては、 μ_1 、 μ_2 ($\mu_1 < \mu_2$) を平均値、 σ_1 、 σ_2 を標準偏差とし、分離信頼度を R として、分離信頼度 R のヒストグラムを、 (μ_1, σ_1) を有する第1の正規分布の確率密度関数 $f_1(R)$ 及び (μ_2, σ_2) を有する第2の正規分布の確率密度関数 $f_2(R)$ でフィッティングすることによって、 μ_1 、 μ_2 、 σ_1 及び σ_2 を推定し、 $f_1(R)$ 、 $f_2(R)$ 、 μ_1 及び μ_2 を使用して前記ソフトマスクを生成する。

【0016】

本実施形態によれば、分離信頼度 R のヒストグラムを正規分布の確率密度関数でフィッティングすることによって、容易にソフトマスクを生成することができる。

40

【0017】

本発明の実施形態による音声認識装置のソフトマスクを生成する方法においては、前記ソフトマスクの値を $S(R)$ 、 $f(R) = f_1(R) + f_2(R)$ として、

$$R < \mu_1 \text{ において } S(R) = 0$$

$$\mu_1 < R < \mu_2 \text{ において } S(R) = f_2(R) / f(R)$$

$$R > \mu_2 \text{ において } S(R) = 1$$

とする。

【0018】

本実施形態によれば、分離信頼度 R のヒストグラムから求めた正規分布の確率密度関数

50

を使用して、容易にソフトマスクを定めることができる。

【0019】

本発明の実施形態による音声認識装置のソフトマスクを生成する方法においては、前記ソフトマスクの値を $S(R)$ 、

$R < \mu_1$ において

【数1】

$$f_1'(R) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

$\mu_1 < R$ において

【数2】

$$f_1'(R) = f_1(R)$$

$R < \mu_2$ において

【数3】

$$f_2'(R) = f_2(R)$$

$\mu_2 < R$ において

【数4】

$$f_2'(R) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

とし、

【数5】

$$f'(R) = f_1'(R) + f_2'(R)$$

として、

【数6】

$$SM(R) = \frac{f_2'(R)}{f'(R)}$$

とする。

【0020】

本実施形態によれば、分離信頼度 R のヒストグラムから求めた正規分布の確率密度関数を使用して、容易にソフトマスクを定めることができる。

【0021】

本発明の実施形態による音声認識装置のソフトマスクを生成する方法においては、 $f_1(R)$ と $f_2(R)$ との交点で

$\mu_1 < R < \mu_2$

を満たす R の値を b とし、

$1 / (1 + \exp(-a(R - b)))$

が

$f_2(R) / f(R)$

とフィッティングするように a を定めて、前記ソフトマスクの値を $S(R)$ として、

$S(R) = 1 / (1 + \exp(-a(R - b)))$

とする。

【0022】

10

20

30

40

50

本実施形態によれば、分離信頼度 R のヒストグラムから求めた正規分布の確率密度関数を使用して、容易にソフトマスクを定めることができる。

【図面の簡単な説明】

【0023】

【図1】本発明の一実施形態による音声認識装置の構成を示す図である。

【図2】音源分離部の構成を示す図である。

【図3】分離信頼度 R の分布を表すヒストグラムである。

【図4】MFMを作成する第1の方法を説明するための図である。

【図5】MFMを作成する第2の方法を説明するための図である。

【図6】MFMを作成する第3の方法を説明するための図である。

【図7】マイクロフォンの位置を示す図である。

【図8】スピーカー及びロボットの配置を示す図である。

【図9】ハードマスクとソフトマスクの概念を示す図である。

【図10】パラメータ探索空間に対する、ソフトマスクの、中央のスピーカーからの単語認識率マップを示す図である。

【図11】ハードマスク及びソフトマスクをベースとする音声認識装置の認識率を示す図である。

【図12】分離信頼度 R の分布を表すヒストグラムを使用した、ソフトMFMの生成方法を示す流れ図である。

【図13】マスクの生成方法を示す流れ図である。

【発明を実施するための形態】

【0024】

図1は、本発明の一実施形態による音声認識装置100の構成を示す図である。音声認識装置100は、音源分離部101、マスク生成部103及び音声認識部105から構成される。

【0025】

音声認識装置100は、複数話者など複数音源の音声を同時認識する。音源分離部101は、たとえば、8チャンネルのマイクロフォンアレイを経て複数音源からの混合音声を受け取る。音源分離部101は、分離音を音声認識部105に送る。また、音源分離部101は、後で説明するように、マスク生成部103が、マスク生成に使用する情報をマスク生成部103に送る。マスク生成部103は、音源分離部101から受け取った情報を使用してマスクを生成し、該マスクを音声認識部105に送る。音声認識部105は、音源分離部101から受け取った分離音の音響特徴量を求め、マスク生成部103から受け取ったマスクを使用して音声認識を行う。音声認識部105、音源分離部101及びマスク生成部103の機能について以下においてさらに説明する。

【0026】

音声認識部

音声認識部105は、ミッシングフィーチャ理論に基づいて、音響特徴量系列及び対応するマスク系列から音素列を出力する。ここで、音響特徴量及びマスクは時間フレームごとに計算される。時間フレームごとに計算された音響特徴量またはマスクを時間に沿って並べたものを系列と呼称する。音声認識部105は、隠れマルコフモデル(HMM)に基づいた認識装置であり、HMMは、従来の自動音声認識システムにおいても普通に使用されている。本実施形態の音声認識部105の自動音声認識方法と、従来の音声認識方法との差異は以下のとおりである。従来の音声認識方法において、最尤パスの推定は、HMMにおける状態遷移及び出力確率に基づいている。この出力確率を推定するプロセスが、本実施形態の音声認識部105において、以下のように修正されている。

【0027】

10

20

30

40

【数 7】

$$M = [M(1), \dots, M(F)]$$

がミッシングフィーチャマスク (MFM) ベクトルであり、

【数 8】

$$M(f)$$

が f 番目の音響特徴量の分離信頼度を表すとする。F は、MFM ベクトルのサイズであり、ある時間フレームの MFM ベクトルは、F 個の要素を含む。

10

【0028】

出力確率

【数 9】

$$b_j(x)$$

は、以下の式で表せる。

【数 10】

$$b_j(x) = \sum_{l=1}^L P(l | S_j) \exp \left\{ \sum_{f=1}^F M(f) \log g(x(f) | l, S_j) \right\} \quad (1)$$

20

但し、 $P(\cdot | \cdot)$ は、確率オペレータである。L は、混合正規分布の混合数を表し、l は、混合正規分布の混合数のインデックスを表す。

【0029】

【数 11】

$$x = [x(1), \dots, x(F)]$$

は、音響特徴量ベクトルであり、F は、音響特徴量ベクトルのサイズである。すなわち、ある時間フレームの音響特徴量ベクトルは、F 個の要素を含む。

【0030】

30

【数 12】

$$S_j$$

は、j 番目の状態であり、

【数 13】

$$g(x(f) | l, S_j)$$

は、j 番目の状態の混合の正規分布である。音響特徴量の分離信頼度の知識が得られなければ、出力確率の式は、従来の式と同じになる。

40

【0031】

音声認識部 105 は、日本語実時間大量単語音声認識エンジンである Julius (参考文献 7) の拡張である Multiband Julius (参考文献 5 及び 6) を使用した。

【0032】

音源分離部

図 2 は、音源分離部 101 の構成を示す図である。図 2 に示すように、音源分離部 101 は、多チャンネルポストフィルタを備えた、幾何学的音源分離 (Geometric Sound Separation, GSS) (参考文献 3、8 及び 11) を使用している。

【0033】

参考文献 9 による GSS アプローチは、確率的な傾きを使用したより速い適応及びより

50

短い時間フレーム推定を提供するように改良されている（参考文献 11）。GSSを使用した最初の分離に、多数音源用のビームフォーマー・ポストフィルタリング（参考文献 11）の一般化に基づくマルチチャンネル・ポストフィルタが続く。このポストフィルタは、最初の分離の間に生成された信号を強化するために、背景ノイズ及び干渉音源の適応スペクトル推定を使用する。

【0034】

音源分離部 101 の音源分離方法の本質的な特徴は、ノイズ推定が定常的な成分と過渡的な成分に分解されていることである。過渡的な成分は、最初の分離段階における出力チャンネル間のリークによると仮定される。

【0035】

このGSS方法は、周波数領域において機能する。

【0036】

【数14】

$$s_m(f, t)$$

が時間フレーム t における離散周波数 f に対する実際の（未知の）音源であるとする。音源

【数15】

$$s_m(f, t)$$

に対応するベクトルは、

【数16】

$$s(f, t)$$

であり、行列

【数17】

$$A(f)$$

は、音源からマイクロフォンへの伝達関数である。マイクロフォンにおいて観察される信号は、以下の式で表現される。

【数18】

$$x(f, t) = A(f) s(f, t) + n(f, t) \quad (2)$$

ここで、

【数19】

$$n(f, t)$$

は、非コヒーレント背景ノイズである。行列

【数20】

$$A(f)$$

は、音源特定アルゴリズムの結果として推定される。全ての伝達関数が単位ゲインを有すると仮定すると、

【数21】

$$A(f)$$

の要素は、以下の式で表現される。

10

20

30

40

50

【 0 0 3 7 】

$$a_{ij}(f) = \exp\{-j 2 \pi f t_{ij}\}$$

(3)

分離結果は、

【 数 2 2 】

$$y(f, t) = W(f, t) x(f, t)$$

と定義され、ここで

【 数 2 3 】

$$W(f, t)$$

10

は、分離行列である。この行列は、参考文献 1 1 に記載された G S S アルゴリズムを使用して推定される。

【 0 0 3 8 】

G S S アルゴリズムの出力は、最初に、参考文献 1 2 によって提案された、最適推定器に基づく周波数領域ポストフィルタによって強化される。

【 0 0 3 9 】

マルチチャネル・ポストフィルタの入力は、G S S の出力

【 数 2 4 】

$$y(f, t) = (y_1(f, t), \dots, y_M(f, t))$$

20

である。マルチチャネル・ポストフィルタの出力

【 数 2 5 】

$$\hat{s}(f, t)$$

は、

【 数 2 6 】

$$\hat{s}(f, t) = G(f, t) y(f, t) \quad (4)$$

と表される。ただし、 $G(f, t)$ は、ゲインである。 $G(f, t)$ の推定値は、スペクトル振幅の最小二乗誤差基準で求める。 $G(f, t)$ を求めるために、ノイズの分散が推定される。

30

【 0 0 4 0 】

ノイズの分散推定値 $\lambda_m(f, t)$ は、

【 数 2 7 】

$$\lambda_m(f, t) = \lambda_m^{stat.}(f, t) + \lambda_m^{leak}(f, t) \quad (5)$$

と表される。ただし、

【 数 2 8 】

$$\lambda_m^{stat.}(f, t)$$

40

と

【 数 2 9 】

$$\lambda_m^{leak}(f, t)$$

は、時間フレーム t における、周波数 f に対する、音源 m のノイズの定常要素の推定値と音源の干渉の推定値である。

【 0 0 4 1 】

50

定常雑音の推定値

【数 3 0】

$$\lambda_m^{stat.}(f, t)$$

は、Minima Controlled Recursive Average (M C R A) (参考文献 1 0) によって求める。

【数 3 1】

$$\lambda_m^{leak}(f, t)$$

は、他の音源からの干渉が、ファクタ によって減少 (典型的には - 1 0 d B - 5 d B) する仮定のもとで、推定される。

【0 0 4 2】

干渉の推定値は、

【数 3 2】

$$\lambda_m^{leak}(f, t) = \eta \sum_{i=0, i \neq m}^{M-1} Z_i(f, t) \quad (6)$$

と表される。ただし、 $Z_i(f, t)$ は、音源 m の平滑化スペクトルで、スペクトル $Y_m(f, t)$ を用いて再帰的に定義される (参考文献 1 1)。

【数 3 3】

$$Z_m(f, t) = \alpha Z_m(f, t - 1) + (1 - \alpha) Y_m(f, t) \quad (7)$$

ただし、 α は - 0 . 7 である。

【0 0 4 3】

マスク生成部

4 8 個の、スペクトルに関連した特徴量の特徴量ベクトルが使用される。ミッシングフイーチャ・マスク (M F M) は、2 4 個の静的スペクトル特徴量及び 2 4 個の動的スペクトル特徴量に対応するベクトルである。ベクトルの各要素は、各特徴量の信頼性を表す。従来の M F M 生成において、2 値の M F M (すなわち、信頼性がある場合は 1 であり、信頼性がない場合は 0 である) が使用されていた。マスク生成部 1 0 3 は、そのベクトルの各要素が 0 . 0 から 1 . 0 の間であるソフト M F M を生成する。ここで、ソフト M F M を生成するとは、ソフト M F M の定義式にしたがって、その値を定めることをいう。

【0 0 4 4】

マスク生成部 1 0 3 は、音源分離部 1 0 1 のマルチチャンネル・ポストフィルタの、入力 $y_m(f, t)$ 、出力

【数 3 4】

$$\hat{s}_m(f, t)$$

及び背景雑音の推定値 $b_n(f, t)$ を使用して M F M を計算する。これらのパラメータは、対象関係伝達関数 (Object related transfer function, ORTF) を使用してマルチチャンネル入力音声から計算される。メル・フィルタバンクを通した変数は、それぞれ、

【数 3 5】

$$Y_m(f, t), \hat{S}_m(f, t), BN(f, t)$$

である。メル・フィルタバンクとは、メル周波数軸上で等間隔に配置されたフィルタ群である。

【0 0 4 5】

分離信頼度 $R(f, t)$ 以下のように定義する。

10

20

30

40

【数 3 6】

$$R(f, t) = \frac{\hat{S}_m(f, t) + BN(f, t)}{Y_m(f, t)} \quad (8)$$

Y_m は、音声

【数 3 7】

$$\hat{S}_m$$

と背景雑音 BN とリークを足し合わせたものからなるため、リークがない場合（他の音源からの混ざりこみがなく、完全に分離できている場合）には分離信頼度が 1 となり、リークが大きくなるにつれて 0 に近い値をとるようになる。

【0 0 4 6】

静的スペクトル特徴量

【数 3 8】

$$[x(1), \dots, x(24)]$$

に対する従来のハード M F M は、以下のように定義される。

【数 3 9】

$$HM_s(f, t) = w_{hard} Q_{hard}(f, t | \theta_{hard}) \quad (9)$$

$$Q_{hard}(f, t | \theta_{hard}) = \begin{cases} 1, & R(f, t) > \theta_{hard} \\ 0, & otherwise \end{cases} \quad (10)$$

ここで、 w_{hard} は、重み係数である。

【数 4 0】

$$0.0 \leq w_{hard} \leq 1.0$$

【0 0 4 7】

動的スペクトル特徴量

【数 4 1】

$$[x(25), \dots, x(48)]$$

に対するハード M F M は、以下のように定義される。

【数 4 2】

$$HM_d(f, t) = \prod_{j=t-2, j \neq t}^{t+2} Q_{hard}(f, j | \theta_{hard}) \quad (11)$$

動的特徴量に対する重み付けされていないハードマスクは、二つの連続するフレーム内の静的特徴量に対するハードマスクが 1 である場合に限り 1 である。

【0 0 4 8】

静的スペクトル特徴量

【数 4 3】

$$[x(1), \dots, x(24)]$$

に対するソフト M F M は、以下のように定義される。

10

20

30

40

【数 4 4】

$$SM_s(f, t) = w Q_{soft}(R(f, t | \theta_{soft}, k)) \quad (12)$$

$$Q_{soft}(x | \theta_{soft}, k) = \begin{cases} \frac{1}{1 + \exp(-k(x - \theta_{soft}))}, & x > \theta_{soft} \\ 0, & otherwise \end{cases} \quad (13)$$

ここで、w は、重み係数である。

10

【0049】

【数 4 5】

$$0.0 \leq w \leq 1.0$$

$$Q_{soft}(\bullet | k, \theta_{soft})$$

は、2 個の調整可能なパラメータを有する修正されたシグモイド関数である。k 及び θ_{soft} は、シグモイド関数の傾きと位置に対応する。シグモイド関数のパラメータの定め方については、後で詳細に説明する。

20

【0050】

動的スペクトル特徴量は、リークノイズ及び静的背景ノイズに対してロバストである。その理由は、隣接する静的スペクトル特徴量の差として定義された動的スペクトル特徴量は、リークノイズ及び静的背景ノイズをキャンセルすることができるからである。静的スペクトル特徴量は、そのようなノイズに対して、動的スペクトル特徴量よりもロバストではない。したがって、動的スペクトル特徴量の寄与が、静的スペクトル特徴量の寄与よりも高い場合には、音声認識率が向上することが期待される。動的スペクトル特徴量の寄与を高くするには、w に小さな値を設定するのが有効である。

【0051】

動的スペクトル特徴量に対するソフト MFM は、以下の式によって定義される。

30

【数 4 6】

$$SM_d(f, t) = \prod_{j=t-2, j \neq t}^{t+2} Q_{soft}(R(f, j | k, \theta_{soft})) \quad (14)$$

【0052】

図 9 は、ハードマスクとソフトマスクの概念を示す図である。図 9 の (a) 及び (c) はハードマスクを示し、図 9 の (b) 及び (d) は、ソフトマスクを示す。図 9 の (a) 及び (b) の横軸は周波数を示し、縦軸はパワーを示す。図 9 の (a) 及び (b) における実線と点線は、それぞれ、クリーンな音声のスペクトル特徴量と歪を受けた音声のスペクトル特徴量を示す。ある周波数における実線と点線との差が歪のパワーを示す。図 9 の (c) 及び (d) の横軸は周波数を示し、縦軸はマスクの値を示す。図 9 の (c) 及び (d) における実線は、マスクの値を示す。図 9 の (c) に示したハードマスクでは、しきい値を使用して歪のある部分のスペクトル特徴量を音声認識における尤度計算から除外する。図 9 の (d) に示したソフトマスクでは、歪のある部分のスペクトル特徴量を歪量に応じて重み付けして尤度計算を行なう。このように、ハードマスクは、歪のある部分のスペクトル特徴量の情報を無駄にしている。したがって、適切に求めたソフトマスクを使用することにより、音声認識率が向上することが期待される。

40

【0053】

上記において、ソフト MFM を、修正されたシグモイド関数を使用して作成した場合に

50

ついて説明した。一般的に、ソフトMFMは、種々の方法によって作成することができる。ここで、ソフトMFMの種々の作成方法について説明する。

【0054】

図12は、分離信頼度Rの分布を表すヒストグラムを使用した、ソフトMFMの生成方法を示す流れ図である。ここで、ソフトマスク(ソフトMFM)を生成するとは、ソフトマスクの定義式を定めることをいう。具体的には、分離信頼度Rの関数としてソフトマスクの定義式を定める。

【0055】

図12のステップS1010において、分離信頼度Rの分布を表すヒストグラムを求める。

10

【0056】

図3は、分離信頼度Rの分布を表すヒストグラムである。横軸は、分離信頼度の値を示し、縦軸は度数を示す。

【0057】

図12のステップS1020において、ステップS1010で求めたヒストグラムに対して、EMアルゴリズム(Expectation-maximization algorithm)を用いて混合正規分布をフィッティングすることにより、第1の正規分布 $f_1(R)$ の平均値及び標準偏差 (μ_1, σ_1) 並びに第2の正規分布 $f_2(R)$ の平均値及び標準偏差 (μ_2, σ_2) を推定する。

【0058】

20

図12のステップS1030において、ステップS1020求めた、 (μ_1, σ_1) 及び (μ_2, σ_2) を使用して以下の方法によりソフトMFMを定めることができる。

【0059】

第1の方法

図4は、MFMを作成する第1の方法を説明するための図である。

【0060】

MFMマスクの値を $S(R)$ 、 $f(R) = f_1(R) + f_2(R)$ として、

$R < \mu_1$ において $S(R) = 0$

$\mu_1 < R < \mu_2$ において $S(R) = f_2(R) / f(R)$

$R > \mu_2$ において $S(R) = 1$

30

とする。

【0061】

第2の方法

図5は、MFMを作成する第2の方法を説明するための図である。

【0062】

MFMマスクの値を $S(R)$ 、

$R < \mu_1$ において

【数47】

$$f_1'(R) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

40

$\mu_1 < R$ において

【数48】

$$f_1'(R) = f_1(R)$$

$R > \mu_2$ において

【数 4 9】

$$f_2'(R) = f_2(R)$$

μ_2 Rにおいて

【数 5 0】

$$f_2'(R) = \frac{1}{\sqrt{2\pi\sigma^2}}$$

10

とし、

【数 5 1】

$$f'(R) = f_1'(R) + f_2'(R)$$

として、

【数 5 2】

$$SM(R) = \frac{f_2'(R)}{f'(R)}$$

20

とする。

【0063】

第3の方法

図6は、MFMを作成する第3の方法を説明するための図である。

【0064】

$f_1(R)$ と $f_2(R)$ との交点で

$$\mu_1 < R < \mu_2$$

を満たすRの値をbとし、

$$1 / (1 + \exp(-a(R - b)))$$

が

$$f_2(R) / f(R)$$

30

とフィッティングするようにaを定めて、MFMマスクの値をS(R)として、

$$S(R) = 1 / (1 + \exp(-a(R - b)))$$

とする。

【0065】

実験

本実施形態による音声認識装置の効率を評価するように、3つの同時音声信号について実験を行った。人間型ロボットに8個の全方位マイクロフォンを取り付けた。マイクロフォンは空中にないので、ロボットの体の伝達関数は、捉えた音に影響を与えた。

【0066】

40

図7は、ロボットに設置されたマイクロフォンの位置を示す図である。図7において、マイクロフォンの位置は矢印で示されている。

【0067】

3個のスピーカーを使用して3つの同時音声信号を生成し、同時音声信号を記録した。反響時間は、0.35秒である。

【0068】

図8は、スピーカー及びロボットの配置を示す図である。1個のスピーカーは、ロボットの正面に配置した。他の2個のスピーカーは、ロボットの左側及び右側の、10、20、30、40、50、60、70、80又は90度の角度に配置した。図8において右側の角度を θ で示し、左側の角度を $-\theta$ で示している。換言すれば、角度 θ を変えながら、

50

9通りの構成で実験を行なった。スピーカーの音量は、全ての場所において同じレベルに設定した。それぞれの構成に対して、3つの異なる単語の200個の組み合わせが実施された。単語は、国際電気通信基礎研究所（ASR）によって配布された、216個の音声的にバランスのとれた単語から選択した。換言すれば、本実施形態による音声認識装置は、各構成において、3つの同時声信号を、200回認識した。

【0069】

式(9)、(12)及び(13)におけるパラメータ θ_{hard} 、 θ_{soft} 、 k 及び w を最適化するように3つの同時音声信号の認識について実験を行った。

【0070】

図13は、マスクの生成方法を示す流れ図である。

10

【0071】

図13のステップS2010において、パラメータを有し、マスクを規定する分離信頼度 R の関数を定める。ハードマスクを規定する関数は、式(9)及び(10)で表され、パラメータは θ_{hard} である。ソフトマスクを規定する関数は、式(12)及び(13)で表され、パラメータは θ_{soft} 、 k 及び w である。

【0072】

図13のステップS2020において、パラメータの探索範囲を定める。

【0073】

表1は、パラメータ探索範囲を示す表である。

20

【表1】

パラメータ	ハードマスク	ソフトマスク
しきい値 θ_{hard}	0.0 - 0.4 (0.05きざみ)	なし
傾き k	なし	80 - 160 (20きざみ)
中央 θ_{soft}	なし	0.0 - 0.4 (0.05きざみ)
重み w	0.0 - 1.0 (0.1きざみ)	0.0 - 1.0 (0.1きざみ)

30

【0074】

図13のステップS2030において、パラメータの探索範囲内でパラメータの値を変化させ、その値を有するマスクを使用した音声認識装置の音声認識率を求める。

【0075】

図13のステップS2040において、音声認識率が最大となるパラメータの値をマスクに使用するパラメータの値とする。

【0076】

結果によれば、ハードマスク θ_{hard} の最適なしきい値(音声認識率を最大とするパラメータ)は、0.1であり、ソフトマスクに設定された最適なパラメータ・セット(音声認識率を最大とするパラメータ・セット)は、

40

$$\{w, \theta_{soft}, k\} = \{0.3, 0.2, 140\}$$

であった。ハードマスク及びソフトマスクに基づいた、中央のスピーカーからの最良の認識率は、それぞれ、93%及び97%であるので、ソフトマスクは、ハードマスクよりもよく機能している。

【0077】

図10は、パラメータ探索空間に対する、ソフトマスクの、中央のスピーカーからの単語認識率マップを示す図である。図10の「しきい値」は、 θ_{soft} を示す。左及び右スピーカーに対しても、マップのピークに設定されるパラメータは、中央のスピーカーに対するマップと同様である。

50

【 0 0 7 8 】

自動音声認識には、Multiband Juliusを使用した。実験においては、分離した単語を認識するのに、三重音音響モデル及び文法ベース言語モデルを使用した。三重音は、3つの状態及び各状態における4つの混合を有するHMMであり、国際電気通信基礎研究所（ASR）によって配布された、216個の音声的にバランスのとれた単語において、訓練される。語彙のサイズは、200語である。

【 0 0 7 9 】

図11は、ハードマスク及びソフトマスクをベースとする音声認識装置の認識率を示す図である。これらの認識率は、全ての探索範囲における、最良の認識率である。横軸は、スピーカーの位置を示し、縦軸は、単語認識率を示す。探索空間の詳細は、表1に示されている。たとえば、横軸上の「30及び左」は、認識目標スピーカーが、中央の30度左側に位置し、他の2個のスピーカーが中央と中央の30度右側に位置することを意味する。横軸上の「60及び中央」は、認識目標スピーカーが、ロボットの正面に位置し、他の2個のスピーカーが中央の60度右側及び左側に位置することを意味する。ソフトマスクをベースとする音声認識装置の語認識率は、ハードマスクをベースとする音声認識装置の語認識率よりも、平均で約5%高い。

10

【 0 0 8 0 】

このように、適切に設計され、調整されたソフトマスクを使用することにより、音声認識装置の、複数音源の音声の同時認識率が向上した。

【 0 0 8 1 】

なお、上記の実施形態においては、分離信頼度Rを使用してソフトマスクを定めた。分離信頼度Rに代えて、音源分離部で求めた入力音声のS/N比（信号・ノイズ比）を使用してソフトマスクの値を設定してもよい。

20

【 0 0 8 2 】

参考文献

[1] Makio Kashino and Tatsuya Hirahara, "One, two, many-judging the number of concurrent talkers," Journal of Acoustic Society of America, vol.99, no.4, pp. Pt .2,2596, 1966.

[2] M. L. Seltzer, B. Raj, and R. M. Stern, "A Bayesian frame work for spectrographic mask estimation for missing feature speech recognition," Speech Communication, vol.43, pp. 379-393, 2004.

30

[3] Shun'ichi Yamamoto, Jean-Marc Valin, Kazuhiro Nakadai, Jean Rouat, Francois Michaud, Tetsuya Ogata, and Hiroshi G. Okuno, "Enhanced Robot Speech Recognition Based on Microphone Array Source Separation and Missing Feature Theory," in Proc. of IEEE CRA-2005, pp. 1489-1494, 2005.

[4] J.Barker, L. Josifovski, M. P. Cooke and P. D. Green, "Soft decision in missing data techniques for robust automatic speech recognition," Proc., ICSLP-2000, 2000.

40

[5] Yoshitaka Nishimura, Takahiro Shinozaki, Koji Iwano, and Sadaoki Furui, "Noise-Robust Speech Recognition Using Multi-Band Spectral Features," in Proc., 148th Acoustical Society of America Meetings, No.1aSC7, 2004.

[6] Multiband Julius, "<http://www.furui.cs.titech.ac.jp/mbandjulius/>".

[7] Tatsuya Kawahara and Akinobu Lee, "Free Software Toolkit for Japanese Large Vocabulary Continuous Speech Recognition," in Proc. of ISCA ICSLP-2000, vol. 4

50

, pp. 476-479, 2000.

[8] Shun'ichi Yamamoto, Kazuhiro Nakadai, Jean-Marc Valin, Jean Rouat, Francois Michaud, Kazunori Komatani, Tetsuya Ogata, and Hiroshi G. Okuno, "Making A Robot Recognize Three Simultaneous Sentences In Real-time," in Proc. of IEEE/RSJIR OS-2005, pp. 897-902, 2005.

[9] Lucas C. Parra and Cristopher V. Alvino, "Geometric Source Separation: Merging Convolutional Source Separation With Geometric Beamforming," IEEE Trans. Speech and Audio Processing, vol. 10, no. 6, pp. 352-362, 2002. 10

[10] Israel Cohen and Baruch Berdugo, "Speech enhancement for non-stationary noise environments," Signal Processing, 81(2), pp. 2403-2418, 2001.

[11] Shun'ichi Yamamoto, Kazuhiro Nakadai, Mikio Nakano, Hiroshi Tsujino, Jean-Marc Valin, Ryu Takeda, Kazunori Komatani, Tetsuya Ogata, and Hiroshi G. Okuno, "Genetic Algorithm-Based Improvement of Robot Hearing Capabilities in Separating and Recognizing Simultaneous Speech Signals," in Proc., IEA/AIE-2006 LNAI4031, 2006, pp. 207-217, Springer-Verlag. 20

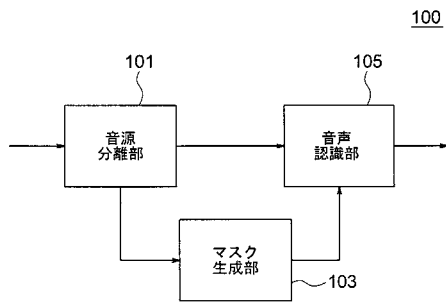
[12] Y. Ephraim and D. Malah, "Speech Enhancement Using Minimum Mean-Square Error Log-Spectral Amplitude Estimator," IEEE Trans. Acoust., Speech, Signal Processing, vol. ASSP-33, no. 2, pp. 443-445, 1985.

【符号の説明】

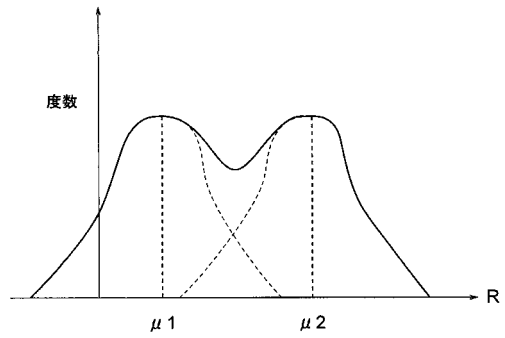
【0083】

100...音声認識装置、101...音源分離部、103...マスク生成部、105...音声認識部

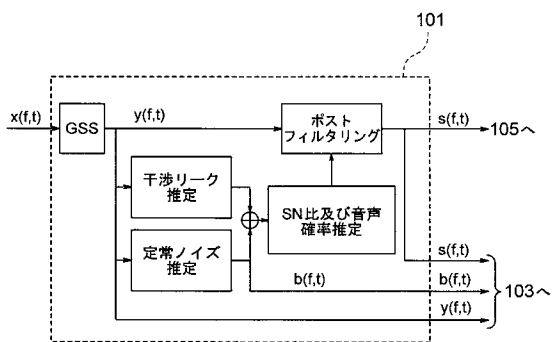
【図1】



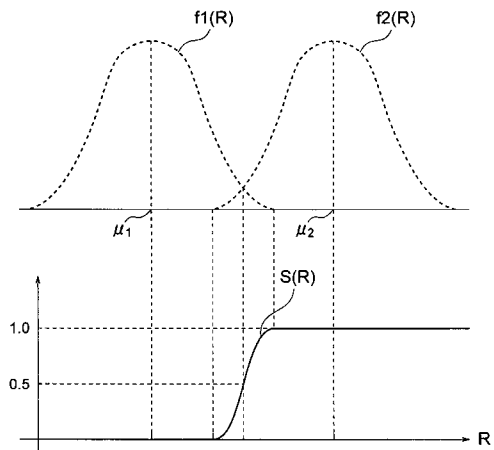
【図3】



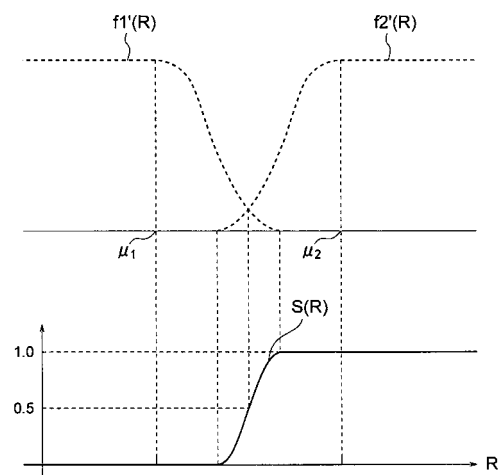
【図2】



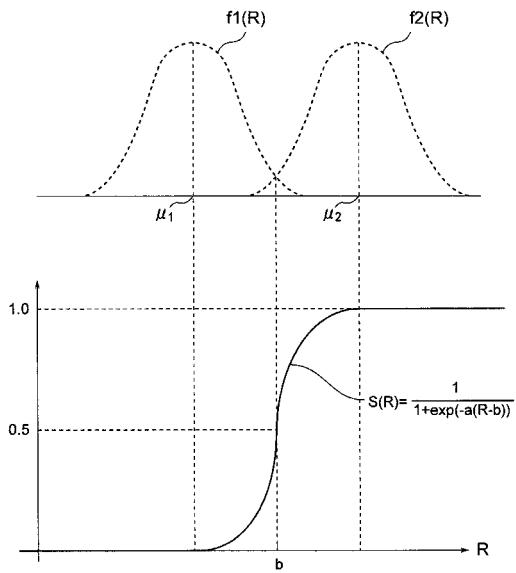
【図4】



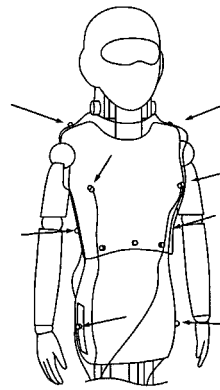
【図5】



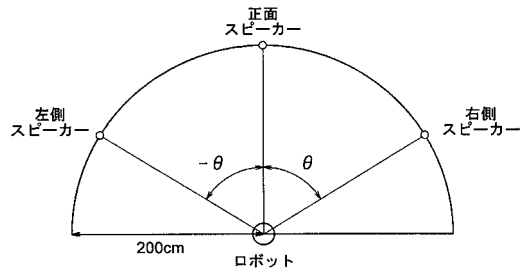
【図6】



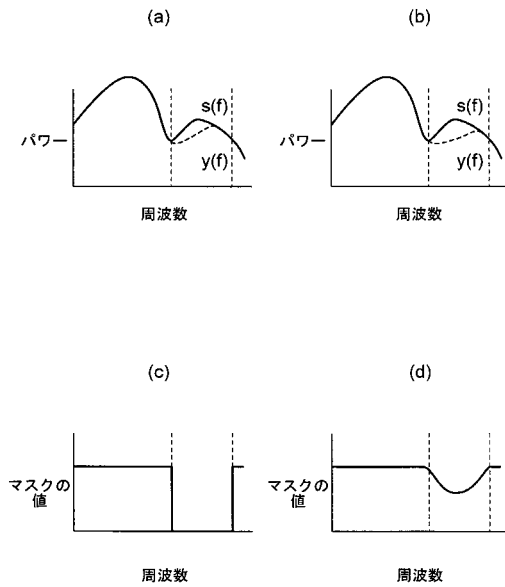
【図7】



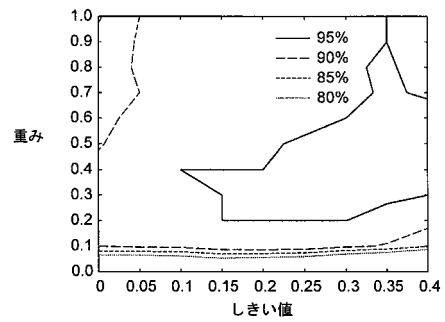
【図8】



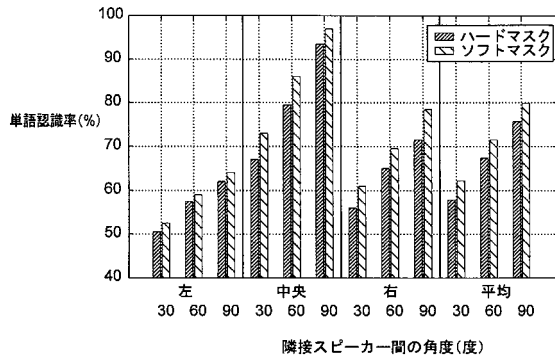
【図9】



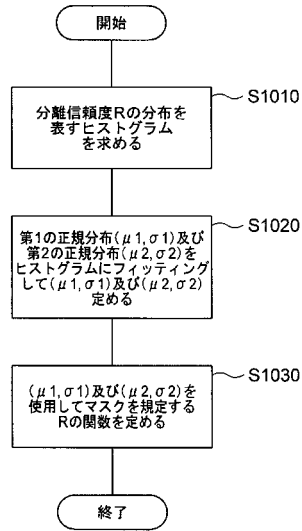
【図10】



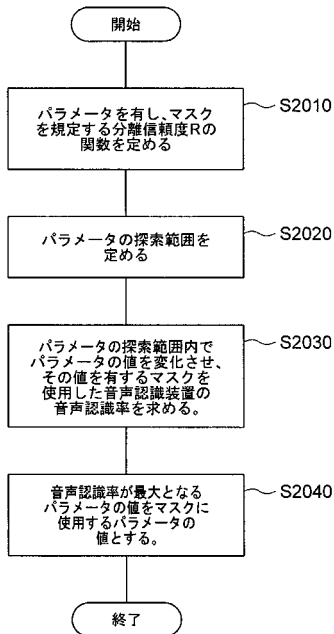
【図11】



【図12】



【図13】



フロントページの続き

審査官 安田 勇太

- (56)参考文献 Yamamoto, S., Valin, J.-M., Nakadai, K., Rouat, J., Michaud, F., Ogata, T., Okuno, H.G., Enhanced Robot Speech Recognition Based on Microphone Array Source Separation and Missing Feature Theory, Proceedings of the 2005 IEEE International Conference on , 2005年 4月18日
武田龍 山本俊一 駒谷和範 尾形哲也 奥乃博, ICAとMFTに基づく音声認識におけるSoft Maskを用いた性能評価, 情報処理学会 第69回(平成19年)全国大会講演論文集(2) 人工知能と認知科学, 日本, 2007年 3月 6日, pp.2-585~2-586

(58)調査した分野(Int.Cl., DB名)

G10L 15/20