

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2005-227511

(P2005-227511A)

(43) 公開日 平成17年8月25日(2005.8.25)

(51) Int. Cl. <sup>7</sup>	F I	テーマコード (参考)
G 1 0 L 15/20	G 1 0 L 3/02 3 0 1 E	5 D 0 1 5
G 1 0 L 11/02	G 1 0 L 3/00 5 1 1	
G 1 0 L 15/02	G 1 0 L 3/00 5 1 3 B	
G 1 0 L 15/04	G 1 0 L 9/08 3 0 1 A	
G 1 0 L 15/28		

審査請求 未請求 請求項の数 15 O L (全 22 頁) 最終頁に続く

(21) 出願番号	特願2004-35618 (P2004-35618)	(71) 出願人	000010076 ヤマハ発動機株式会社 静岡県磐田市新貝2500番地
(22) 出願日	平成16年2月12日 (2004. 2. 12)	(74) 代理人	100066980 弁理士 森 哲也
		(74) 代理人	100075579 弁理士 内藤 嘉昭
		(74) 代理人	100103850 弁理士 崔 秀▲てつ▼
		(72) 発明者	有宗 伸泰 静岡県磐田市新貝2500番地 ヤマハ発動機株式会社内
		(72) 発明者	赤坂 貴志 静岡県磐田市新貝2500番地 ヤマハ発動機株式会社内
		Fターム(参考)	5D015 DD02 DD03 EE04

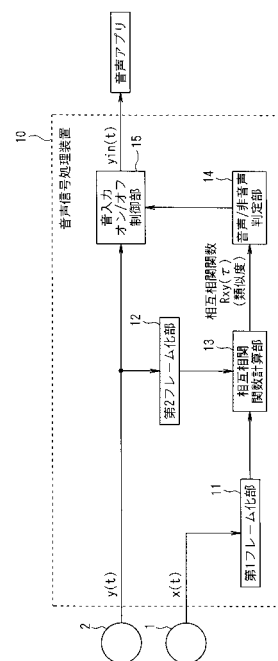
(54) 【発明の名称】 対象音検出方法、音信号処理装置、音声認識装置及びプログラム

(57) 【要約】

【課題】マイクの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系の構築を可能にする。

【解決手段】音声信号処理装置は、発話音又は雑音のいずれか一方を受音するように配置されている単一指向性マイク1と、発話音及び雑音を受音するように配置されている無指向性マイク2と、単一指向性マイク1に入力された音声信号  $x(t)$  と無指向性マイク2に入力された音声信号  $y(t)$  との相互相関関数  $R_{xy}(\tau)$  を算出するためのフレーム化部11、12及び相互相関関数計算部13と、相互相関関数計算部13が算出した相互相関関数  $R_{xy}(\tau)$  に基づいて、発話音の発話区間を検出する音声/非音声判定部14とを備える。

【選択図】 図3



**【特許請求の範囲】****【請求項 1】**

無指向性マイクで検出対象音及び雑音を受音し、単一指向性マイクで前記検出対象音又は前記雑音のいずれか一方を受音し、前記無指向性マイクに入力された声信号と単一指向性マイクに入力された音信号とを比較し、その比較結果に基づいて、前記検出対象音を検出することを特徴とする対象音検出方法。

**【請求項 2】**

前記無指向性マイクに入力された音信号と単一指向性マイクに入力された音信号との比較により相関度を得て、その相関度に基づいて、前記検出対象音を検出することを特徴とする請求項 1 記載の対象音検出方法。

10

**【請求項 3】**

前記無指向性マイクに入力された音信号のパワースペクトルと、単一指向性マイクに入力された音信号のパワースペクトルとを比較して、その比較結果に基づいて、前記検出対象音を検出することを特徴とする請求項 1 又は 2 記載の対象音検出方法。

**【請求項 4】**

前記無指向性マイクに入力された音信号と単一指向性マイクに入力された音信号との比較により得た相関度、及び前記無指向性マイクに入力された音信号のパワースペクトルと、単一指向性マイクに入力された音信号のパワースペクトルとの比較結果に基づいて、前記検出対象音を検出することを特徴とする請求項 1 乃至 3 のいずれか 1 に記載の対象音検出方法。

20

**【請求項 5】**

検出対象音及び雑音を受音するように配置されている無指向性マイクと、  
前記検出対象音又は前記雑音のいずれか一方を受音するように配置されている単一指向性マイクと、  
前記無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較する比較手段と、  
前記比較手段の比較結果に基づいて、前記検出対象音を検出する対象音検出手段と、  
を備えることを特徴とする音信号処理装置。

**【請求項 6】**

前記比較手段は、前記無指向性マイクに入力された音信号と前記単一指向性マイクに入力された音信号との相関度を算出し、前記対象音検出手段は、前記比較手段が算出した相関度と所定の第 1 しきい値とを比較して、前記検出対象音を検出することを特徴とする請求項 5 記載の音信号処理装置。

30

**【請求項 7】**

前記比較手段は、前記無指向性マイク及び単一指向性マイクに入力された各音信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した 2 つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、

前記対象音検出手段は、前記パワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第 2 しきい値とを比較して、前記検出対象音を検出することを特徴とする請求項 5 又は 6 記載の音信号処理装置。

40

**【請求項 8】**

前記比較手段は、前記無指向性マイク及び単一指向性マイクに入力された各音信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した 2 つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、

前記比較手段は、前記無指向性マイクに入力された音信号と前記単一指向性マイクに入力された音信号との相関度を算出し、前記対象音検出手段は、前記比較手段が算出した相関度と所定の第 1 しきい値との比較結果と、前記比較手段のパワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第 2 しきい値との比較結果とに基づいて、

50

前記検出対象音を検出することを特徴とする請求項 5 乃至 7 のいずれか 1 に記載の音信号処理装置。

【請求項 9】

前記無指向性マイクに入力された音信号及び単一指向性マイクに入力された音信号を時分割してフレーム化するフレーム化手段を備えており、

前記比較手段は、前記フレーム化手段から出力されるフレーム単位で、前記無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較し、前記対象音検出手段は、前記比較手段の比較結果に基づいて、前記フレーム化手段から出力されるフレーム単位で前記検出対象音を検出することを特徴とする請求項 5 乃至 8 のいずれか 1 に記載の音信号処理装置。

10

【請求項 10】

発話音及び雑音を受音するように配置されている無指向性マイクと、

前記発話音又は前記雑音のいずれか一方を受音するように配置されている単一指向性マイクと、

前記無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較する比較手段と、

前記比較手段の比較結果に基づいて、前記発話音の発話区間を検出する発話区間検出手段と、

前記発話区間検出手段が検出した前記発話音の発話区間について、音声認識処理を行う音声認識処理手段と、

20

を備えることを特徴とする音声認識装置。

【請求項 11】

前記比較手段は、前記無指向性マイクに入力された音信号と前記単一指向性マイクに入力された音信号との相関度を算出し、前記発話区間検出手段は、前記比較手段が算出した相関度と所定の第 1 しきい値とを比較して、前記発話音の発話区間を検出することを特徴とする請求項 10 記載の音声認識装置。

【請求項 12】

前記比較手段は、前記無指向性マイク及び単一指向性マイクに入力された各音信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した 2 つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、

30

前記発話区間検出手段は、前記パワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第 2 しきい値とを比較して、前記発話音の発話区間を検出することを特徴とする請求項 10 又は 11 記載の音声認識装置。

【請求項 13】

前記比較手段は、前記無指向性マイク及び単一指向性マイクに入力された各音信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した 2 つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、

前記比較手段は、前記無指向性マイクに入力された音信号と前記単一指向性マイクに入力された音信号との相関度を算出し、前記発話区間検出手段は、前記比較手段が算出した相関度と所定の第 1 しきい値との比較結果と、前記比較手段のパワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第 2 しきい値との比較結果とに基づいて、前記発話音の発話区間を検出することを特徴とする請求項 10 乃至 12 のいずれか 1 に記載の音声認識装置。

40

【請求項 14】

前記無指向性マイクに入力された音信号及び単一指向性マイクに入力された音信号を時分割してフレーム化するフレーム化手段を備えており、

前記比較手段は、前記フレーム化手段から出力されるフレーム単位で、前記無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較し、前記発話

50

区間検出手段は、前記比較手段の比較結果に基づいて、前記フレーム化手段から出力されるフレーム単位で前記発話音の発話区間を検出し、前記音声認識処理手段は、前記発話区間検出手段が検出したフレーム単位の前記発話音の発話区間について、音声認識処理を行うことを特徴とする請求項10乃至13のいずれか1に記載の音声認識装置。

【請求項15】

無指向性マイクで受音した検出対象音及び雑音の音信号と単一指向性マイクで受音した前記検出対象音又は前記雑音のいずれか一方の音信号とを比較し、その比較結果に基づいて、前記検出対象音を検出する処理をコンピュータに実行させることを特徴とするプログラム。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明は、入力音中から検出対象音を検出する対象音検出方法及びこれを実現するプログラム、入力された音信号を処理する音信号処理装置、並びに入力された発話音について音声認識処理を行う音声認識装置に関する。

【背景技術】

【0002】

音声は、人間の用いる種々の通信の形態の中でも最も根源的なものであると同時に、他のどの情報送出方法よりも高速度に情報を送り出すことのできる優れた通信手段である。このようなことから、音声は、古くから現在に至るまで人間の通信手段の根幹を担っている。

20

音声認識技術は、そのような音声を認識するための技術である。音声認識とは、その音声に含まれる情報の中で、最も基本的な意味内容に関する情報、つまり音韻情報をコンピュータなどにより抽出し、その抽出内容を判定することである。近年では、計算機プロセッサ技術の飛躍的な発達と、インターネットに代表される高度な情報ネットワークの構築により、様々な分野においてマン・マシンインタフェースとしての音声認識技術の適用が試みられている。

現在の音声認識システムの認識性能は、確率・統計的手法により格段に向上しており、理想的な環境下での音声や接話マイクロホンで収録された近距離音声などでは、非常に高い認識率が得られるようになっている。

30

【発明の開示】

【発明が解決しようとする課題】

【0003】

実環境下の音声認識は、学習データと観測データとの間の環境や発話内容のミスマッチ等により、その認識率が劣化する。また、受音系となる接話マイクヘッドセットの装着によりユーザが受ける負担や不快感は大きく、音声認識システム実用化の大きな障害のひとつになっている。

また、S/N比の低下や背景雑音、室内残響の影響などにより認識が困難な遠隔音声に関し、複数の遠隔マイクロホンを用いた音声認識手法の研究が多くなされている。その代表的なものとして、マイクロホンアレーを用いる手法が挙げられる。この手法では、音源位置検出処理、目的音強調処理、雑音抑制処理、の3つの空間的な信号処理を行なうことができる。このような手法により遠隔音声の音声認識が盛んに研究されている。

40

【0004】

しかし、この手法は、正確な話者方向同定処理のために複数のマイクロホンを一定間隔にて固定配置する必要があるため、小型化、携帯化が困難であるため、様々な環境・状況下での音声入力への応用が難しく、用途が限定されるという問題がある。

本発明は、前述の問題に鑑みてなされたものであり、マイクロホンの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系の構築を可能にする対象音検出方法、音信号処理装置、音声認識装置及びプログラムの提供を目的とする。

50

## 【課題を解決するための手段】

## 【0005】

請求項1記載の対象音検出方法は、無指向性マイクで検出対象音及び雑音を受音し、単一指向性マイクで前記検出対象音又は前記雑音のいずれか一方を受音し、前記無指向性マイクに入力された音声信号と単一指向性マイクに入力された音声信号とを比較し、その比較結果に基づいて、前記検出対象音を検出することを特徴とする。

また、請求項2記載の対象音検出方法は、請求項1記載の対象音検出方法において、前記無指向性マイクに入力された音声信号と単一指向性マイクに入力された音声信号との比較により相関度を得て、その相関度に基づいて、前記検出対象音を検出することを特徴とする。

10

## 【0006】

また、請求項3記載の対象音検出方法は、請求項1又は2記載の対象音検出方法において、前記無指向性マイクに入力された音声信号のパワースペクトルと、単一指向性マイクに入力された音声信号のパワースペクトルとを比較して、その比較結果に基づいて、前記検出対象音を検出することを特徴とする。

また、請求項4記載の対象音検出方法は、請求項1乃至3のいずれか1に記載の対象音検出方法において、前記無指向性マイクに入力された音声信号と単一指向性マイクに入力された音声信号との比較により得た相関度、及び前記無指向性マイクに入力された音声信号のパワースペクトルと単一指向性マイクに入力された音声信号のパワースペクトルとの比較結果に基づいて、前記検出対象音を検出することを特徴とする。

20

## 【0007】

また、請求項5記載の音声信号処理装置は、検出対象音及び雑音を受音するように配置されている無指向性マイクと、前記検出対象音又は前記雑音のいずれか一方を受音するように配置されている単一指向性マイクと、前記無指向性マイクに入力された音声信号と、単一指向性マイクに入力された音声信号とを比較する比較手段と、前記比較手段の比較結果に基づいて、前記検出対象音を検出する対象音検出手段と、を備えることを特徴とする。

## 【0008】

また、請求項6記載の音声信号処理装置は、請求項5記載の音声信号処理装置において、前記比較手段が、前記無指向性マイクに入力された音声信号と前記単一指向性マイクに入力された音声信号との相関度を算出し、前記対象音検出手段が、前記比較手段が算出した相関度と所定の第1しきい値とを比較して、前記検出対象音を検出することを特徴とする。

30

また、請求項7記載の音声信号処理装置は、請求項5又は6記載の音声信号処理装置において、前記比較手段が、前記無指向性マイク及び単一指向性マイクに入力された各音声信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した2つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、前記対象音検出手段が、前記パワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第2しきい値とを比較して、前記検出対象音を検出することを特徴とする。

## 【0009】

また、請求項8記載の音声信号処理装置は、請求項5乃至7のいずれか1に記載の音声信号処理装置において、前記比較手段が、前記無指向性マイク及び単一指向性マイクに入力された各音声信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した2つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、前記比較手段が、前記無指向性マイクに入力された音声信号と前記単一指向性マイクに入力された音声信号との相関度を算出し、前記対象音検出手段は、前記比較手段が算出した相関度と所定の第1しきい値との比較結果と、前記比較手段のパワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第2しきい値との比較結果とに基づいて、前記検出対象音を検出することを特徴とする。

40

## 【0010】

また、請求項9記載の音声信号処理装置は、請求項5乃至8のいずれか1に記載の音声信号

50

処理装置において、前記無指向性マイクに入力された音信号及び単一指向性マイクに入力された音信号を時分割してフレーム化するフレーム化手段を備えており、前記比較手段が、前記フレーム化手段から出力されるフレーム単位で、前記無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較し、前記対象音検出手段は、前記比較手段の比較結果に基づいて、前記フレーム化手段から出力されるフレーム単位で前記検出対象音を検出することを特徴とする。

【0011】

また、請求項10記載の音声認識装置は、発話音及び雑音を受音するように配置されている無指向性マイクと、前記発話音又は前記雑音のいずれか一方を受音するように配置されている単一指向性マイクと、前記無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較する比較手段と、前記比較手段の比較結果に基づいて、前記発話音の発話区間を検出する発話区間検出手段と、前記発話区間検出手段が検出した前記発話音の発話区間について、音声認識処理を行う音声認識処理手段と、を備えることを特徴とする。

10

【0012】

また、請求項11記載の音声認識装置は、請求項10記載の音声認識装置において、前記比較手段が、前記無指向性マイクに入力された音信号と前記単一指向性マイクに入力された音信号との相関度を算出し、前記発話区間検出手段が、前記比較手段が算出した相関度と所定の第1しきい値とを比較して、前記発話音の発話区間を検出することを特徴とする。

20

【0013】

また、請求項12記載の音声認識装置は、請求項10又は11記載の音声認識装置において、前記比較手段が、前記無指向性マイク及び単一指向性マイクに入力された各音信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した2つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、前記発話区間検出手段が、前記パワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第2しきい値とを比較して、前記発話音の発話区間を検出することを特徴とする。

【0014】

また、請求項13記載の音声認識装置は、請求項10乃至12のいずれか1に記載の音声認識装置において、前記比較手段が、前記無指向性マイク及び単一指向性マイクに入力された各音信号のパワースペクトルをそれぞれ算出するパワースペクトル算出手段と、前記パワースペクトル算出手段が算出した2つのパワースペクトルの比を算出するパワースペクトル比算出手段とを備えており、前記比較手段が、前記無指向性マイクに入力された音信号と前記単一指向性マイクに入力された音信号との相関度を算出し、前記発話区間検出手段は、前記比較手段が算出した相関度と所定の第1しきい値との比較結果と、前記比較手段のパワースペクトル比算出手段が算出した前記パワースペクトルの比と所定の第2しきい値との比較結果とに基づいて、前記発話音の発話区間を検出することを特徴とする。

30

【0015】

また、請求項14記載の音声認識装置は、請求項10乃至13のいずれか1に記載の音声認識装置において、前記無指向性マイクに入力された音信号及び単一指向性マイクに入力された音信号を時分割してフレーム化するフレーム化手段を備えており、前記比較手段が、前記フレーム化手段から出力されるフレーム単位で、前記無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較し、前記発話区間検出手段が、前記比較手段の比較結果に基づいて、前記フレーム化手段から出力されるフレーム単位で前記発話音の発話区間を検出し、前記音声認識処理手段は、前記発話区間検出手段が検出したフレーム単位の前記発話音の発話区間について、音声認識処理を行うことを特徴とする。

40

【0016】

50

また、請求項 15 記載のプログラムは、無指向性マイクで受音した検出対象音及び雑音の音信号と単一指向性マイクで受音した前記検出対象音又は前記雑音のいずれか一方の音信号とを比較し、その比較結果に基づいて、前記検出対象音を検出する処理をコンピュータに実行させることを特徴とする。

無指向性マイクで検出対象音及び雑音を受音し、単一指向性マイクで検出対象音又は雑音のいずれか一方を受音するようにした場合において、無指向性マイク及び単一指向性マイクが配置されている環境で雑音だけを発したときと検出対象音及び雑音を発したときとで、無指向性マイク及び単一指向性マイクに入力された音信号が異なってくる。なお、検出対象音には、人間が発する発話音の他、物体が発する音も含まれる。

【0017】

そこで、請求項 1、5、15 記載の発明では、無指向性マイクに入力された音信号と単一指向性マイクに入力された音信号とを比較することで、検出対象音又は発話音を検出している。また、請求項 10 記載の発明では、検出対象音が人間が発する発話音であり、検出対象音の検出として、発話音の音声区間の検出を行っている。

また、請求項 2、6、11 記載の発明では、無指向性マイクに入力された音信号と単一指向性マイクに入力された音信号との相関度により、検出対象音の検出又は発話音の発話区間の検出を行っている。

【0018】

また、請求項 3、7、12 記載の発明では、無指向性マイクに入力された音信号のパワースペクトルと単一指向性マイクに入力された音信号のパワースペクトルとを比較することで、検出対象音の検出又は発話音の発話区間の検出を行っている。

また、請求項 4、8、13 記載の発明では、無指向性マイクに入力された音信号と単一指向性マイクに入力された音信号との相関度と、無指向性マイクに入力された音信号のパワースペクトルと単一指向性マイクに入力された音信号のパワースペクトルとの比較結果とに基づいて、検出対象音の検出又は発話音の発話区間の検出を行っている。

また、請求項 9、14 記載の発明では、前記無指向性マイクに入力された音信号及び単一指向性マイクに入力された音信号を時分割してフレーム化し、フレーム単位でその後の処理を行う。

【発明の効果】

【0019】

本発明によれば、無指向性マイクで検出対象音（又は発話音）及び雑音を受音し、単一指向性マイクで前記検出対象音（又は発話音）又は前記雑音のいずれか一方を受音するように、無指向性マイク及び単一指向性マイクを配置する限り、検出対象音（又は発話音の音声区間）を検出することができる。これにより、マイクロホンの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系の構築が可能になる。

【発明を実施するための最良の形態】

【0020】

本発明を実施するための最良の形態（以下、実施形態という。）を図面を参照しながら詳細に説明する。

第 1 の実施形態は、図 1 に示すように、第 1 及び第 2 マイク 1, 2 に入力された音声信号を処理する音声信号処理装置 10 である。

第 1 マイク 1 は単一指向性マイクであり、第 2 マイク 2 は無指向性マイクであり、第 1 及び第 2 マイク 1, 2 は例えば装着型マイクである。第 1 及び第 2 マイク 1, 2 は、図 2 に示すように、第 1 及び第 2 マイク 1, 2 をできるだけ近づけて配置するとともに、単一指向性マイクである第 1 マイク 1 をその指向方向が音源（ユーザ）の位置に対して反対側となるように配置している。また、第 1 マイク 1 の指向方向に、雑音源が存在している。なお、図 2 に示す点線は、雑音源を基準にした第 1 マイク 1 の指向特性を示し、図 2 に示す一点鎖線は、第 2 マイク 2 の指向特性を示す。

このように第 1 及び第 2 マイク 1, 2 を配置すると、雑音源からの音は、第 1 及び第 2

10

20

30

40

50

マイク 1, 2 で受音でき、音源 (ユーザ) からの音は第 2 マイク 2 だけが受音できるようになる。

【0021】

図 3 は、第 1 の実施形態の音声信号処理装置 10 の構成を示す。

図 3 に示すように、音声信号処理装置 10 は、第 1 及び第 2 フレーム化部 11, 12、相互相関関数計算部 13、音声 / 非音声判定部 14 並びに音入力オン / オフ制御部 15 を備えている。

第 1 及び第 2 マイク 1, 2 から入力された 2 ch の音声信号はそれぞれ、第 1 及び第 2 フレーム化部 11, 12 に入力される。また、第 2 マイク 2 から入力された音声信号は、音入力オン / オフ制御部 15 に入力される。ここで、第 1 マイク 1 に入力された音声信号を  $x(t)$  とし、第 2 マイク 2 に入力された音声信号を  $y(t)$  とする。

10

【0022】

第 1 フレーム化部 11 では、第 1 マイク 1 から入力された音声信号  $x(t)$  を時分割でフレーム化 (或いはフレーム分割) して、複数フレームにした音声信号  $x(t)$  を相互相関関数計算部 13 に出力する。また、第 2 フレーム化部 12 では、第 2 マイク 2 から入力される音声信号  $y(t)$  を時分割でフレーム化 (或いはフレーム分割) して、複数フレームにした音声信号  $y(t)$  を相互相関関数計算部 13 に出力する。ここで、第 1 及び第 2 フレーム化部 11, 12 は、入力されてくる音声信号  $x(t)$ ,  $y(t)$  を所定時間間隔でサンプリングしていき、所定のサンプル数を 1 フレームとして次々にフレーム化していく。

20

【0023】

相互相関関数計算部 13 は、第 1 フレーム化部 11 から出力されるフレームと、第 2 フレーム化部 12 から出力されるフレームとを比較する。すなわち、第 1 マイク 1 に入力された音声信号  $x(t)$  と、第 2 マイク 2 に入力された音声信号  $y(t)$  とをフレーム単位で比較する。その比較結果として、相互相関関数計算部 13 は、下記 (1) 式により、相互相関関数  $R_{xy}(\tau)$  を算出する。

【0024】

【数 1】

$$R_{xy}(\tau) = \frac{1}{T} \int_{-T/2}^{T/2} x(t)y(t-\tau)dt \quad \cdots (1)$$

30

【0025】

ここで、 $\tau$  は第 1 マイク 1 と第 2 マイク 2 との間の距離によって決まる遅延時間である。また、 $T$  はフレーム長である。

前述したように第 1 及び第 2 マイク 1, 2 をできるだけ近づけて配置している場合には、遅延時間  $\tau$  を近似的に 0 とおくことができる。しかし、後述するような本発明の効果を満たす限り、第 1 マイク 1 と第 2 マイク 2 とを離して配置することは可能であり、この場合、遅延時間  $\tau$  を適切に与える必要がある。すなわち例えば、第 1 マイク 1 と第 2 マイク 2 との間の距離を 10 cm にしている場合には、その 10 cm 相当分の遅延時間  $\tau$  を与えて、相互相関関数  $R_{xy}(\tau)$  を算出する。このようにすれば、第 1 マイク 1 と第 2 マイク 2 との間の距離を考慮して、相互相関関数  $R_{xy}(\tau)$  を得ることができ、精度よく相互相関関数  $R_{xy}(\tau)$  を得ることができる。

40

【0026】

このように算出された相互相関関数  $R_{xy}(\tau)$  はフレーム単位で各音声信号  $x(t)$ ,  $y(t)$  の波形形状の類似度を示す値となる。具体的には、相関関係を求める 2 つの音声信号  $x(t)$ ,  $y(t)$  が似ているほど、相互相関関数  $R_{xy}(\tau)$  は大きい値となり、相関関係を求める 2 つの音声信号  $x(t)$ ,  $y(t)$  が異なっているほど、相互相関関数  $R_{xy}(\tau)$  は 0 に近くなる。相互相関関数計算部 13 は、このような相互相関関数  $R_{xy}(\tau)$  を音声 / 非音声判定部 14 に出力する。

50



## 【0027】

音声／非音声判定部14は、相互相関関数 $R \times y$  ( )に基づいて音声区間(発話区間)と非音声区間(非発話区間)とを判定する。具体的には、次のように音声区間と非音声区間とを判定する。

前述したように、音源(ユーザ)と雑音源に対して図2のように第1及び第2マイク1, 2を配置することで、雑音源からの音を第1及び第2マイク1, 2で受信し、音源(ユーザ)からの音を第2マイク2だけで受信している。

## 【0028】

一方、相互相関関数 $R \times y$  ( )は、前述したように、相関関係を求める2つの音声信号 $x(t)$ ,  $y(t)$ が似ているほど大きい値となり、相関関係を求める2つの音声信号 $x(t)$ ,  $y(t)$ が異なっているほど0に近くなる。

このようなことから、雑音源からの音だけを第1及び第2マイク1, 2で受信している場合には、同じ音声信号が第1及び第2マイク1, 2に入力されているので、すなわち、第1及び第2マイク1, 2の入力音声信号のS/N比が同程度になるので、相互相関関数 $R \times y$  ( )は大きい値になる。一方、音源(ユーザ)から発話があった場合には、その発話を第2マイク2だけが受信するので、第1及び第2マイク1, 2それぞれに異なる音声信号が入力されるようになり、すなわち第2マイク2の入力音声信号のS/N比の方が大きくなるので、相互相関関数 $R \times y$  ( )は0に向かって減少する。

## 【0029】

このように、音源(ユーザ)から発話があった場合には相互相関関数 $R \times y$  ( )は0に向かって減少することから、音声／非音声判定部14は、相互相関関数 $R \times y$  ( )と判定用しきい値(類似度を示すしきい値) $r_1$ とを比較して、音声区間を判定する。すなわち、音声／非音声判定部14は、相互相関関数 $R \times y$  ( )が判定用しきい値 $r_1$ 未満の場合( $R \times y$  ( ) <  $r_1$ )、音声区間と判定し、それ以外の場合( $R \times y$  ( )  $\geq r_1$ )、非音声区間と判定する。ここで、判定用しきい値 $r_1$ は例えば実験により得る。そして、音声／非音声判定部14は、このような判定をフレーム単位で行う。音声／非音声判定部14は、その判定結果をフレーム単位で音入力オン／オフ制御部15に出力する。

## 【0030】

音入力オン／オフ制御部15には、第2マイク2からの音声信号 $y(t)$ が入力されており、音入力オン／オフ制御部15は、音声／非音声判定部14の判定結果に基づいて、第2マイク2からの音声信号 $y(t)$ の後段への出力のオンとオフとを切り換える。具体的には、音声／非音声判定部14が音声区間と判定した場合、音入力オン／オフ制御部15は、オン制御として当該音声区間に対応する音声信号 $y(t)$ の区間を後段に出力して、音声／非音声判定部14が非音声区間と判定した場合、音入力オン／オフ制御部15は、オフ制御として当該非音声区間に対応する音声信号 $y(t)$ の区間を後段に出力しないようにする。

## 【0031】

以上のように音声信号処理装置10が構成されている。この音声信号処理装置10における一連の動作は次のようになる。

まず、第1及び第2フレーム化部11, 12が、第1及び第2マイク1, 2から入力された2chの音声信号 $x(t)$ ,  $y(t)$ をそれぞれフレーム化し、フレーム単位で音声信号 $x(t)$ ,  $y(t)$ を相互相関関数計算部13に出力する。

## 【0032】

相互相関関数計算部13では、第1及び第2フレーム化部11, 12それぞれから出力されるフレーム単位の音声信号 $x(t)$ ,  $y(t)$ について相互相関関数 $R \times y$  ( )を算出して、算出した相互相関関数 $R \times y$  ( )を音声／非音声判定部14に出力する。

音声／非音声判定部14では、相互相関関数 $R \times y$  ( )と判定用しきい値 $r_1$ とを比較し、相互相関関数 $R \times y$  ( )に対応するフレームが音声区間のものか、非音声区間のものかを判定する。そして、音声／非音声判定部14は、その判定結果を音入力オン／

10

20

30

40

50

オフ制御部 15 に出力する。

【0033】

音入力オン/オフ制御部 15 は、音声/非音声判定部 14 が音声区間と判定した場合、オン制御として第 2 マイク 2 からの音声信号  $y(t)$  を後段に出力して、音声/非音声判定部 14 が非音声区間と判定した場合、オフ制御として第 2 マイク 2 からの音声信号  $y(t)$  を後段に出力しないようにする。このとき、音入力オン/オフ制御部 15 から出力される音声信号  $y(t)$  は、音源(ユーザ)からの音と雑音源からの音とからなる音声信号となる。

【0034】

このように、音声信号処理装置 10 は、第 2 マイク 2 への入力音中の発話区間(音声区間)を検出することができる。 10

例えば、第 1 マイク 1, 2 と音声アプリケーションとの間にこのような音声信号処理装置 10 を備えることで、音声アプリケーションは、確実に発話区間についての処理を行うことができる。ここで、音声アプリケーションとしては、音声認識システム、放送システム、携帯電話及びトランシーバが挙げられる。例えば、音声アプリケーションが音声認識システムであるとすれば、音声認識システムは、音声信号処理装置 10 が出力する発話区間の音声信号に基づいて音声認識できるようになる。

【0035】

次に第 1 の実施形態における効果を説明する。

前述したように、無指向性マイクである第 2 のマイク 2 で発話音及び雑音を受音し、単一指向性マイクである第 1 マイク 1 で雑音を受音し、第 1 マイク 1 で受音した雑音の音声信号と第 2 マイク 2 で受音した発話音及び雑音からなる音声信号との比較により相関度を得て、その相関度に基づいて、発話音の発話区間を特定している。 20

【0036】

これにより、第 2 のマイク 2 で発話音及び雑音を受音し、かつ第 1 マイク 1 で雑音を受音するように第 1 及び第 2 のマイク 1, 2 を配置する限り、マイクの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系を構築することができる。

また、精度よく発話区間を検出することができる。そして、このように音声信号処理装置 10 が検出した発話区間の音声信号を利用することにより、音声認識システムでは、高認識率、低誤認識率の音声認識が可能になり、また、携帯電話やトランシーバでは、信頼性の高いハンズフリー半二重通信が可能になり、放送システムでは、通信システムの送信電力低減が可能になる。 30

【0037】

次に第 2 の実施形態を説明する。

この第 2 の実施形態も前述の第 1 の実施形態と同様、前記図 1 に示したように 2 つのマイク 1, 2 に入力された音声信号を処理する音声信号処理装置 10 である。そして、音声信号処理装置 10 の構成は、前述の第 1 の実施形態と同様、図 3 に示したような構成になる。しかし、第 2 の実施形態では、第 1 及び第 2 マイク 1, 2 の配置を前述の第 1 の実施形態における配置と異ならせている。 40

前述の第 1 の実施形態では、第 1 マイク 1 に単一指向性マイクを用い、第 2 マイク 2 に無指向性マイクを用い、前記図 2 に示したように、雑音源からの音を第 1 及び第 2 マイク 1, 2 で受音し、音源(ユーザ)からの音を第 1 マイク 1 だけで受音するように、第 1 及び第 2 マイク 1, 2 を配置している。

【0038】

一方、この第 2 の実施形態では、音源(ユーザ)からの音を第 1 及び第 2 マイク 1, 2 で受音し、雑音源からの音を第 1 マイク 1 だけで受音するようにしている。具体的には、第 1 マイク 1 に無指向性マイクを用い、第 2 マイク 2 に単一指向性マイクを用いる。そして、図 4 に示すように、第 1 及び第 2 マイク 1, 2 をできるだけ近づけて配置するとともに、単一指向性マイクである第 2 マイク 2 を、その指向方向が音源(ユーザ)に向かい、 50

かつその指向方向外に雑音源が位置されるように、配置する。なお、図4に示す点線は、第1マイク1の指向特定を示し、図4に示す一点鎖線は、音源（ユーザ）を基準にした第2マイク2の指向特性を示す。

【0039】

このように第1及び第2マイク1, 2を配置した場合、前述の第1の実施形態と比較し、特に相互相関関数計算部13で算出される相互相関関数 $R \times y$ （ ）が異なる傾向を示すようになる。

すなわち、音源（ユーザ）からの音を第1及び第2マイク1, 2で受音し、雑音源からの音を第2マイク2だけが受音しているため、雑音源からの音だけを第1マイク1で受音している場合には、第1及び第2マイク1, 2それぞれに異なる音声信号が入力されるようになり、このとき相互相関関数 $R \times y$ （ ）は0に近い値になる。一方、音源（ユーザ）から発話があった場合には、その発話を第1及び第2マイク1, 2で受音するので、ほぼ同じ音声信号が第1及び第2マイク1, 2に入力されるようになり、これにより、相互相関関数 $R \times y$ （ ）は大きい値になる。このとき、第2マイク2の入力音声信号のS/N比は高くなり、第1マイク1の入力音声信号のS/N比は、第2マイク2ほどではないが、高くなる。

【0040】

このように、音源（ユーザ）から発話があった場合には、相互相関関数 $R \times y$ （ ）が大きくなり、第2の実施形態で得る相互相関関数 $R \times y$ （ ）は、前述の第1の実施形態とは反対の傾向を示すようになる。

このようなことから、第2の実施形態では、音声/非音声判定部14は、相互相関関数 $R \times y$ （ ）と判定用しきい値（類似度を示すしきい値） $r_2$ とを比較して、相互相関関数 $R \times y$ （ ）が判定用しきい値 $r_2$ より大きい場合（ $R \times y$ （ ） $> r_2$ ）、音声区間と判定し、それ以外の場合（ $R \times y$ （ ） $< r_2$ ）、非音声区間と判定する。ここで、判定用しきい値 $r_2$ は例えば実験により得る。

【0041】

そして、前述の第1の実施形態と同様に、音声/非音声判定部14が音声区間と判定した場合、音入力オン/オフ制御部15は、オン制御として第2マイク2からの音声信号 $y(t)$ を後段に出力して、音声/非音声判定部34が非音声区間と判定した場合、音入力オン/オフ制御部15は、オフ制御として第2マイク2からの音声信号 $y(t)$ を後段に出力しないようにする。このとき、音入力オン/オフ制御部15から出力される音声信号 $y(t)$ は、音源（ユーザ）からの音のみからなる音声信号となる。

このように、第2の実施形態の音声信号処理装置10は、第2マイク2への入力音中の発話区間（音声区間）を検出することができる。

【0042】

次に第2の実施形態における効果を説明する。

前述したように、無指向性マイクである第1のマイク1で発話音及び雑音を受音し、単一指向性マイクである第2マイク2で発話音を受音し、第1マイク1で受音した発話音及び雑音からなる音声信号と第2マイク2で受音した発話音の音声信号との比較により相関度を得て、その相関度に基づいて、発話音の発話区間を特定している。

【0043】

これにより、第1のマイク1で発話音及び雑音を受音し、かつ第2マイク2で発話音を受音するように第1及び第2のマイク1, 2を配置する限り、マイクの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系を構築することができる。

また、精度よく発話区間を検出することができる。そして、このように音声信号処理装置10が検出した発話区間の音声信号を利用することにより、音声認識システムでは、高認識率、低誤認識率の音声認識が可能になり、また、携帯電話やトランシーバでは、信頼性の高いハンズフリー半二重通信が可能になり、放送システムでは、通信システムの送信電力低減が可能になる。

10

20

30

40

50

## 【 0 0 4 4 】

次に第 3 の実施形態を説明する。

この第 3 の実施形態も前述の第 1 の実施形態と同様、前記図 1 に示したように 2 つのマイク 1, 2 に入力された音声信号を処理する音声信号処理装置 10 である。

前述の第 1 及び第 2 の実施形態では、相互相関関数計算部 13 により音声信号  $x(t)$ ,  $y(t)$  の相互相関関数  $R_{xy}(\tau)$  を算出し、この相互相関関数  $R_{xy}(\tau)$  に基づいて音声区間と非音声区間とを判定している。

## 【 0 0 4 5 】

これに対して、第 3 の実施形態の音声信号処理装置 10 は、音声信号  $x(t)$ ,  $y(t)$  それぞれのパワースペクトルを算出し、このパワースペクトルに基づいて音声区間と非音声区間とを判定するように構成されている。図 5 は、その第 3 の実施形態の音声信号処理装置 10 の構成を示す。

図 5 に示すように、音声信号処理装置 10 は、第 1 及び第 2 フレーム化部 11, 12、第 1 及び第 2 パワースペクトラム計算部 21, 22、パワー比計算部 23、音声 / 非音声判定部 24 並びに音入力オン / オフ制御部 15 を備えている。各部の処理内容は次のようになる。

## 【 0 0 4 6 】

なお、第 1 及び第 2 フレーム化部 11, 12 及び音入力オン / オフ制御部 15 については、前述の第 1 の実施形態のものと同様な処理を行うので、その説明を省略する。また、前述の第 1 の実施形態と同様に、第 1 マイク 1 は単一指向性マイクであり、第 2 マイク 2 は無指向性マイクである。さらに、第 1 及び第 2 マイク 1, 2 の配置についても、前記図 2 に示したような配置にしている。これにより、雑音源からの音を第 1 及び第 2 マイク 1, 2 で受音し、音源 (ユーザ) からの音を第 2 マイク 2 だけで受音している。

## 【 0 0 4 7 】

第 1 及び第 2 フレーム化部 11, 12 それぞれで複数フレームにされた音声信号  $x(t)$ ,  $y(t)$  は第 1 及び第 2 パワースペクトラム計算部 21, 22 に入力される。

第 1 パワースペクトラム計算部 21 は、フレーム単位で音声信号  $x(t)$  の第 1 パワースペクトル値  $P_x(\omega)$  を算出し、その算出した第 1 パワースペクトル値  $P_x(\omega)$  をパワー比計算部 23 に出力する。また、第 2 パワースペクトラム計算部 22 は、フレーム単位で音声信号  $y(t)$  の第 2 パワースペクトル値  $P_y(\omega)$  を算出し、その算出した第 2

パワースペクトル値  $P_y(\omega)$  をパワー比計算部 23 に出力する。  
 パワー比計算部 23 は、下記 (2) 式により、第 1 パワースペクトラム計算部 21 からの第 1 パワースペクトル値  $P_x(\omega)$  と、第 2 パワースペクトラム計算部 22 からの第 2 パワースペクトル値  $P_y(\omega)$  との比 (以下、パワー比という。)  $P_{xy}(\omega)$  を算出する。

## 【 0 0 4 8 】

## 【 数 2 】

$$P_{xy}(\omega) = G_{xy} \frac{P_x(\omega)}{P_y(\omega)} \quad \dots (2)$$

## 【 0 0 4 9 】

ここで、 $G_{xy}$  は、第 1 及び第 2 マイク 1, 2 の感度によって決まる補正係数である。このように算出されたパワー比  $P_{xy}(\omega)$  はフレーム単位で各音声信号  $x(t)$ ,  $y(t)$  の波形形状の類似度を示す値となる。パワー比計算部 23 は、このようなパワー比  $P_{xy}(\omega)$  を音声 / 非音声判定部 24 に出力する。

音声 / 非音声判定部 24 は、パワー比  $P_{xy}(\omega)$  に基づいて音声区間と非音声区間とを判定する。具体的には、次のように音声区間と非音声区間とを判定する。

前述したように、音源 (ユーザ) と雑音源に対して前記図 2 のように第 1 及び第 2 マイク 1, 2 を配置することで、雑音源からの音を第 1 及び第 2 マイク 1, 2 で受音し、音源

10

20

30

40

50

(ユーザ)からの音を第2マイク2だけで受音している。

【0050】

これにより、雑音源からの音だけを第1及び第2マイク1, 2で受音している場合には、同じ音声信号が第1及び第2マイク1, 2に入力されているので、すなわち第1及び第2マイク1, 2の受音感度が同程度であるので、このときに第1及び第2パワースペクトラム計算部21, 22で算出される第1及び第2パワースペクトル値 $P_x(\quad)$ ,  $P_y(\quad)$ は同程度になる。一方、音源(ユーザ)から発話があった場合には、その発話を第2マイク2だけが受音するので、すなわち第2マイク2の受音感度の方が大きくなるので、このときに第1パワースペクトル値 $P_x(\quad)$ よりも第2パワースペクトル値 $P_y(\quad)$ の方が大きくなる。このとき、パワー比計算部23が算出するパワー比 $P_x y(\quad)$ は小さくなる。 10

【0051】

なお、このとき、雑音源や音源(ユーザ)の特性に応じて、所定の周波数域のパワースペクトル値 $P_x(\quad)$ ,  $P_y(\quad)$ が特に変化する。

このように、音源(ユーザ)から発話があった場合にはパワー比 $P_x y(\quad)$ は小さくなることから、音声/非音声判定部24は、パワー比 $P_x y(\quad)$ と判定用しきい値(類似度を示すしきい値)  $p_1$ とを比較して、音声区間を判定する。

【0052】

ここで、第1及び第2パワースペクトラム計算部21, 22では、パワースペクトル値 $P_x(\quad)$ ,  $P_y(\quad)$ を所定の周波数域を対象として得ている。よって、パワー比 $P_x y(\quad)$ は、各周波数帯について得ることができる。 20

このようなことから、パワースペクトル値 $P_x(\quad)$ ,  $P_y(\quad)$ について各周波数で得ているパワー比 $P_x y(\quad)$ の総和平均値を算出し、判定では、その総和平均値と判定用しきい値  $p_1$ とを比較する。ここで、判定用しきい値  $p_1$ は例えば実験により得る。

【0053】

なお、判定対象としてパワースペクトル値 $P_x(\quad)$ ,  $P_y(\quad)$ の全周波数域の総和平均値を用いることに限定されるものではない。例えば、音源(ユーザ)の特性を示す特定の周波数帯のパワー比 $P_x y(\quad)$ の総和平均値と判定用しきい値  $p_1$ とを比較したり、雑音源の特性を示す特定の周波数帯のパワー比 $P_x y(\quad)$ の平均値と判定用しきい値  $p_1$ とを比較したり、又は音源(ユーザ)の特性を示す特定の周波数帯のパワー比 $P_x y(\quad)$ と雑音源の特性を示す特定の周波数帯のパワー比 $P_x y(\quad)$ との平均値と判定用しきい値  $p_1$ とを比較したりしてもよい。この場合、それに応じて、判定用しきい値  $p_1$ を設定する。 30

【0054】

そして、音声/非音声判定部24は、パワー比 $P_x y(\quad)$ が判定用しきい値  $p_1$ 未満の場合( $P_x y(\quad) < p_1$ )、音声区間と判定し、それ以外の場合( $P_x y(\quad) \geq p_1$ )、非音声区間と判定する。ここで、音声/非音声判定部24は、このような判定をフレーム単位で行う。そして、音声/非音声判定部24は、その判定結果を音入力オン/オフ制御部15に出力する。 40

【0055】

以上のように第3の実施形態の音声信号処理装置10が構成されている。この音声信号処理装置10における一連の動作は次のようになる。

まず、第1及び第2フレーム化部11, 12が、第1及び第2マイク1, 2から入力された2chの音声信号 $x(t)$ ,  $y(t)$ をそれぞれフレーム化し、フレーム単位で音声信号 $x(t)$ ,  $y(t)$ を第1及び第2パワースペクトラム計算部21, 22に出力する。

【0056】

パワースペクトラム計算部21, 22ではそれぞれ、第1及び第2フレーム化部11, 12それぞれから出力されるフレーム単位の音声信号 $x(t)$ ,  $y(t)$ について第1及 50

び第2パワースペクトル値  $P_x(\quad)$  ,  $P_y(\quad)$  を算出して、算出した第1及び第2パワースペクトル値  $P_x(\quad)$  ,  $P_y(\quad)$  をパワー比計算部23に出力する。

パワー比計算部23では、パワースペクトラム計算部21, 22それぞれから出力される第1及び第2パワースペクトル値  $P_x(\quad)$  ,  $P_y(\quad)$  について、フレーム単位でパワー比  $P_x/y(\quad)$  を算出して、算出したパワー比  $P_x/y(\quad)$  を音声/非音声判定部24に出力する。

【0057】

音声/非音声判定部24では、パワー比  $P_x/y(\quad)$  と判定用しきい値  $p_1$  とを比較し、パワー比  $P_x/y(\quad)$  に対応するフレームが音声区間のものか、非音声区間のものかを判定する。そして、音声/非音声判定部24は、その判定結果を音入力オン/オフ制御部15に出力する。

10

音入力オン/オフ制御部15では、第2マイク2からの音声信号  $y(t)$  の後段への出力のオンとオフとを切り換える。具体的には、音声/非音声判定部24が音声区間と判定した場合、音入力オン/オフ制御部15は、オン制御として第2マイク2からの音声信号  $y(t)$  を後段に出力して、音声/非音声判定部24が非音声区間と判定した場合、音入力オン/オフ制御部15は、オフ制御として第2マイク2からの音声信号  $y(t)$  を後段に出力しないようにする。このとき、音入力オン/オフ制御部15から出力される音声信号  $y(t)$  は、音源(ユーザ)からの音と雑音源からの音とからなる音声信号となる。

このように、第3の実施形態の音声信号処理装置10は、第2マイク2への入力音中の発話区間(音声区間)を検出することができる。

20

【0058】

次に第3の実施形態における効果を説明する。

前述したように、無指向性マイクである第2のマイク2で発話音及び雑音を受音し、単一指向性マイクである第1マイク1で雑音を受音し、第2マイク2で受音した発話音及び雑音からなる音声信号のパワースペクトルと、第1マイク1で受音した雑音の音声信号のパワースペクトルとを比較して、その比較結果に基づいて、前記発話音の発話区間を特定している。

【0059】

これにより、第2のマイク2で発話音及び雑音を受音し、かつ第1マイク1で雑音を受音するように第1及び第2のマイク1, 2を配置する限り、マイクの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系を構築することができる。

30

また、精度よく発話区間を検出することができる。そして、このように音声信号処理装置10が検出した発話区間の音声信号を利用することにより、音声認識システムでは、高認識率、低誤認識率の音声認識が可能になり、また、携帯電話やトランシーバでは、信頼性の高いハンズフリー半二重通信が可能になり、放送システムでは、通信システムの送信電力低減が可能になる。

【0060】

次に第4の実施形態を説明する。

この第4の実施形態も前述の第3の実施形態と同様、音声信号  $x(t)$  ,  $y(t)$  それぞれのパワースペクトルを算出し、このパワースペクトルに基づいて音声区間と非音声区間とを判定する音声信号処理装置10である。音声信号処理装置10の構成は、前述の第3の実施形態と同様、前記図5に示したような構成になる。そして、第4の実施形態では、第1及び第2マイク1, 2の配置を前述の第3の実施形態における配置と異ならせている。

40

【0061】

前述の第3の実施形態では、前述の第1の実施形態と同様、第1マイク1に単一指向性マイクを用い、第2マイク2に無指向性マイクを用い、前記図2に示したように、雑音源からの音を第1及び第2マイク1, 2で受音し、音源(ユーザ)からの音を第2マイク2だけで受音するように、第1及び第2マイク1, 2を配置している。

50

一方、この第4の実施形態では、前述の第2の実施形態と同様、音源（ユーザ）からの音を第1及び第2マイク1, 2で受音し、雑音源からの音を第1マイク1だけで受音している。具体的には、第1マイク1に無指向性マイクを用い、第2マイク2に単一指向性マイクを用いる。そして、前記図4に示したように、第1及び第2マイク1, 2をできるだけ近づけて配置するとともに、単一指向性マイクである第2マイク2を、その指向方向が音源（ユーザ）に向かい、かつその指向方向外に雑音源が位置されるように、配置する。

【0062】

このように第1及び第2マイク1, 2を配置した場合、パワー比計算部23で算出されるパワー比 $P \times y$  ( )は次のような傾向を示す。

音源（ユーザ）からの音を第1及び第2マイク1, 2で受音し、雑音源からの音を第1マイク1だけで受音するようにしているため、雑音源からの音だけを第1マイク1で受音している場合には、第1マイク1の受音感度の方が大きくなるので、第1パワースペクトル値 $P \times$  ( )が大きくなり、これにより、パワー比 $P \times y$  ( )が大きくなる。一方、音源（ユーザ）から発話があった場合には、その発話を第1及び第2マイク1, 2で受音するので、これにより、第2パワースペクトル値 $P y$  ( )も大きくなり、パワー比 $P \times y$  ( )が小さくなる。

【0063】

このように、音源（ユーザ）から発話があった場合には、パワー比 $P \times y$  ( )が小さくなる傾向を示すようになる。

このようなことから、第4の実施形態では、音声/非音声判定部24は、パワー比 $P \times y$  ( )と判定用しきい値（類似度を示すしきい値） $p^2$ とを比較して、パワー比 $P \times y$  ( )が判定用しきい値 $p^2$ 未満の場合（ $P \times y$  ( ) <  $p^2$ ）、音声区間と判定し、それ以外の場合（ $P \times y$  ( )  $\geq p^2$ ）、非音声区間と判定する。ここで、判定用しきい値 $p^2$ は例えば実験により得る。

【0064】

そして、前述の第3の実施形態と同様に、音入力オン/オフ制御部15は、音声/非音声判定部24が音声区間と判定した場合、オン制御として第2マイク2からの音声信号 $y(t)$ を後段に出力して、音声/非音声判定部24が非音声区間と判定した場合、オフ制御として第2マイク2からの音声信号 $y(t)$ を後段に出力しないようにする。このとき、音入力オン/オフ制御部15から出力される音声信号 $y(t)$ は、音源（ユーザ）からの音のみからなる音声信号となる。

【0065】

このように、第4の実施形態の音声信号処理装置10は、第2マイク2への入力音中の発話区間（音声区間）を検出することができる。

次に第4の実施形態における効果を説明する。

前述したように、無指向性マイクである第1のマイク1で発話音及び雑音を受音し、単一指向性マイクである第2マイク2で発話音を受音し、第1マイク1で受音した発話音及び雑音からなる音声信号のパワースペクトルと、第2マイク2で受音した発話音の音声信号のパワースペクトルとを比較して、その比較結果に基づいて、前記発話音の発話区間を特定している。

【0066】

これにより、第1のマイク1で発話音及び雑音を受音し、かつ第2マイク2で発話音を受音するように第1及び第2のマイク1, 2を配置する限り、マイクの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系を構築することができる。

また、精度よく発話区間を検出することができる。そして、このように音声信号処理装置10が検出した発話区間の音声信号を利用することにより、音声認識システムでは、高認識率、低誤認識率の音声認識が可能になり、また、携帯電話やトランシーバでは、信頼性の高いハンズフリー半二重通信が可能になり、放送システムでは、通信システムの送信電力低減が可能になる。

10

20

30

40

50

## 【0067】

次に第5の実施形態を説明する。

前述の第1及び第2の実施形態では、相互相関関数計算部13により音声信号 $x(t)$ 、 $y(t)$ の相互相関関数 $R_{xy}(\tau)$ を算出し、この相互相関関数 $R_{xy}(\tau)$ に基づいて音声区間と非音声区間とを判定し、また、前述の第3及び第4の実施形態では、音声信号 $x(t)$ 、 $y(t)$ それぞれのパワースペクトル $P_x(\omega)$ 、 $P_y(\omega)$ を算出し、このパワースペクトル $P_x(\omega)$ 、 $P_y(\omega)$ (具体的にはパワー比 $P_{xy}(\omega)$ )に基づいて音声区間と非音声区間とを判定している。第5の実施形態では、第1の実施形態(第2の実施形態)の処理と、第3の実施形態(第4の実施形態)の処理とを組み合わせた処理により、音声区間と非音声区間とを判定している。すなわち、第5の実施形態では、音声信号 $x(t)$ 、 $y(t)$ の相互相関関数 $R_{xy}(\tau)$ を算出するとともに、音声信号 $x(t)$ 、 $y(t)$ それぞれのパワースペクトル $P_x(\omega)$ 、 $P_y(\omega)$ を算出し、相互相関関数 $R_{xy}(\tau)$ とパワースペクトル $P_x(\omega)$ 、 $P_y(\omega)$ (具体的にはパワー比 $P_{xy}(\omega)$ )との両面から音声区間と非音声区間とを判定している。図6は、それを実現する音声信号処理装置10の構成を示す。

10

## 【0068】

図6に示すように、音声信号処理装置10は、第1及び第2フレーム化部11、12、相互相関関数計算部13、音声/非音声判定部14、音入力オン/オフ制御部15、第1及び第2パワースペクトラム計算部21、22、パワー比計算部23、並びに音声/非音声判定部31を備えている。

20

このような構成において、第1及び第2フレーム化部11、12、相互相関関数計算部13、音声/非音声判定部14、音入力オン/オフ制御部15、第1及び第2パワースペクトラム計算部21、22、並びにパワー比計算部23は、前述の第1及び第2の実施形態と同様な処理を行う。

## 【0069】

すなわち、第1及び第2フレーム化部11、12は、第1及び第2マイク1、2から入力された2chの音声信号 $x(t)$ 、 $y(t)$ をそれぞれフレーム化し、フレーム単位で音声信号 $x(t)$ 、 $y(t)$ を相互相関関数計算部13に出力する。また、第1及び第2フレーム化部11、12はそれぞれ、フレーム単位で各音声信号 $x(t)$ 、 $y(t)$ を第1及び第2パワースペクトラム計算部21、22それぞれに出力する。

30

## 【0070】

相互相関関数計算部13は、第1及び第2フレーム化部11、12それぞれから出力されるフレーム単位の音声信号 $x(t)$ 、 $y(t)$ について相互相関関数 $R_{xy}(\tau)$ を算出して、算出した相互相関関数 $R_{xy}(\tau)$ を音声/非音声判定部31に出力する。

また、第1及び第2パワースペクトラム計算部21、22は、第1及び第2フレーム化部11、12それぞれから出力されるフレーム単位の音声信号 $x(t)$ 、 $y(t)$ について第1及び第2パワースペクトル値 $P_x(\omega)$ 、 $P_y(\omega)$ を算出して、算出した第1及び第2パワースペクトル値 $P_x(\omega)$ 、 $P_y(\omega)$ をパワー比計算部23に出力する。パワー比計算部23は、パワースペクトラム計算部21、22それぞれから出力される第1及び第2パワースペクトル値 $P_x(\omega)$ 、 $P_y(\omega)$ について、フレーム単位でパワー比 $P_{xy}(\omega)$ を算出して、算出したパワー比 $P_{xy}(\omega)$ を音声/非音声判定部31に出力する。

40

## 【0071】

音声/非音声判定部31では、次のような判定処理を行う。この音声/非音声判定部31で行う処理が第5の実施形態において特有の処理になる。ここで、前記図2に示したように、雑音源からの音を第1及び第2マイク1、2で受音し、音源(ユーザ)からの音を第2マイク2だけで受音するように、第1及び第2マイク1、2を配置した場合(第1又は第3の実施形態の場合)と、前記図4に示したように、雑音源からの音を第1マイク1だけで受音し、音源(ユーザ)からの音を第1及び第2マイク2で受音するように、第1及び第2マイク1、2を配置した場合(第2又は第4の実施形態の場合)とを分けして説

50



明する。

【0072】

先ず、前記図2に示したように、雑音源からの音を第1及び第2マイク1, 2で受音し、音源(ユーザ)からの音を第2マイク2だけで受音するように、第1及び第2マイク1, 2を配置した場合について説明する。

この場合、音声/非音声判定部31は、音源(ユーザ)から発話があった場合に相互相関関数  $R \times y(\quad)$  が0に向かって減少することから、相互相関関数  $R \times y(\quad)$  が判定用しきい値  $r_1$  未満の場合 ( $R \times y(\quad) < r_1$ )、音声区間とし、それ以外の場合 ( $R \times y(\quad) \geq r_1$ )、非音声区間とする第1判定結果を得る。また、音声/非音声判定部31は、音源(ユーザ)から発話があった場合にはパワー比  $P \times y(\quad)$  が小さくなることから、パワー比  $P \times y(\quad)$  が判定用しきい値  $p_1$  未満の場合 ( $P \times y(\quad) < p_1$ )、音声区間とし、それ以外の場合 ( $P \times y(\quad) \geq p_1$ )、非音声区間とする第2判定結果を得る。

10

【0073】

そして、音声/非音声判定部31は、前記第1及び第2判定結果に基づいて、音声区間の最終的な判定結果を得る。例えば、音声/非音声判定部31は、第1及び第2判定結果が共に音声区間である判定結果となった場合、最終的な判定結果を音声区間とする。または、音声/非音声判定部31は、第1判定結果又は第2判定結果の少なくとも一方が音声区間である判定結果となった場合、最終的な判定結果を音声区間とする。そして、音声/非音声判定部31は、それ以外の場合、最終的な判定結果を非音声区間とする。

20

【0074】

一方、前記図4に示したように、雑音源からの音を第1マイク1だけで受音し、音源(ユーザ)からの音を第1及び第2マイク2で受音するように、第1及び第2マイク1, 2を配置した場合には、次のような判定を行う。

音声/非音声判定部31は、音源(ユーザ)から発話があった場合に相互相関関数  $R \times y(\quad)$  が大きくなることから、相互相関関数  $R \times y(\quad)$  が判定用しきい値  $r_2$  より大きい場合 ( $R \times y(\quad) \geq r_2$ )、音声区間とし、それ以外の場合 ( $R \times y(\quad) < r_2$ )、非音声区間とする第1判定結果を得る。また、音声/非音声判定部31は、音源(ユーザ)からの音の出力(発話)があった場合にはパワー比  $P \times y(\quad)$  が小さくなることから、パワー比  $P \times y(\quad)$  が判定用しきい値  $p_2$  未満の場合 ( $P \times y(\quad) < p_2$ )、音声区間とし、それ以外の場合 ( $P \times y(\quad) \geq p_2$ )、非音声区間とする第2判定結果を得る。

30

【0075】

そして、音声/非音声判定部31は、前記第1及び第2判定結果に基づいて、音声区間の最終的な判定結果を得る。例えば、音声/非音声判定部31は、第1及び第2判定結果が共に音声区間である判定結果となった場合、最終的な判定結果を音声区間とする。または、音声/非音声判定部31は、第1判定結果又は第2判定結果の少なくとも一方が音声区間である判定結果となった場合、最終的な判定結果を音声区間とする。そして、音声/非音声判定部31は、それ以外の場合、最終的な判定結果を非音声区間とする。

【0076】

以上のようにして音声/非音声判定部31は、音声区間又は非音声区間を判定している。そして、音声/非音声判定部31は、その判定結果を音入力オン/オフ制御部15に出力する。

40

音入力オン/オフ制御部15は、音声/非音声判定部31が音声区間と判定した場合、オン制御として第2マイク2からの音声信号  $y(t)$  を後段に出力して、音声/非音声判定部31が非音声区間と判定した場合、オフ制御として第2マイク2からの音声信号  $y(t)$  を後段に出力しないようにする。このとき、第1及び第2のマイク1, 2の配置が前記図2に示した配置であれば、音入力オン/オフ制御部15から出力される音声信号  $y(t)$  は、音源(ユーザ)からの音と雑音源からの音とからなる音声信号となる。また、第1及び第2のマイク1, 2の配置が前記図4に示した配置であれば、音入力オン/オフ制

50

御部 15 から出力される音声信号  $y(t)$  は、音源（ユーザ）からの音のみからなる音声信号となる。

【0077】

次に第5の実施形態における効果を説明する。

前述したように、無指向性マイクで発話音及び雑音を受音し、単一指向性マイクで発話音又は雑音のいずれか一方を受音し、無指向性マイクで受音した発話音及び雑音からなる音声信号と単一指向性マイクで受音した発話音又は雑音のいずれか一方の音声信号の比較により相関度を得ている。その一方で、無指向性マイクで受音した発話音及び雑音からなる音声信号のパワースペクトルと、単一指向性マイクで受音した発話音又は雑音のいずれか一方の音声信号のパワースペクトルとを比較して、その比較結果としてパワー比を得ている。そして、前記相関度とパワー比との両方に基づいて、最終的に前記発話音の発話区間を特定している。

10

【0078】

このように、相関度とパワー比との両方に基づいて最終的に前記発話音の発話区間を特定することで、その特定を精度よく行うことができる。そして、このように音声信号処理装置 10 が検出した発話区間の音声信号を利用することにより、音声認識システムでは、高認識率、低誤認識率の音声認識が可能になり、また、携帯電話やトランシーバでは、信頼性の高いハンズフリー半二重通信が可能になり、放送システムでは、通信システムの送信電力低減が可能になる。

【0079】

また、前述の第1乃至第4の実施形態と同様に、無指向性マイクで発話音及び雑音を受音し、かつ単一指向性マイクで発話音又は雑音のいずれか一方を受音するように第1及び第2のマイク 1, 2 を配置する限り、マイクの取り付け位置の変化等による環境の変化、話者の移動や姿勢の変化等による音源の移動に対してロバストな受音系を構築することができる。

20

【0080】

なお、前述の実施形態では、第1及び第2マイク 1, 2 から入力された音声信号  $x_1(t)$ ,  $x_2(t)$  を、直接第1及び第2フレーム化部 11, 12 にそれぞれ入力しているが、具体的には、第1及び第2マイク 1, 2 から入力された音声信号  $x_1(t)$ ,  $x_2(t)$  を、A/D（アナログ/デジタル）変換した後、第1及び第2フレーム化部 11, 12 に入力するようにする。また、前述の実施形態では、第2マイク 2 に入力された音声信号  $x_1(t)$  を、音入力オン/オフ制御部 15 に入力しているが、第2マイク 2 に入力された、フレーム化した音声信号  $x_1(t)$  を音入力オン/オフ制御部 15 に入力する。これらの仕様を、例えば前述の第1の実施形態の音声信号処理装置 10 の構成に適用すると、図7に示すような構成になる。

30

【0081】

この図7に示すように、第1及び第2マイク 1, 2 から入力された音声信号  $x_1(t)$ ,  $x_2(t)$  をそれぞれ、第1及び第2 A/D変換部 41, 42 で A/D変換した後、第1及び第2フレーム化部 11, 12 に入力する。また、第2 A/D変換部 32 で A/D変換された信号は、第2フレーム化部 12 でフレーム化されてから音入力オン/オフ制御部 15 に入力される。ここで、第1及び第2 A/D変換部 41, 42 で A/D変換されたデータ形式は、例えば 11025 Hz、16 bit、リニア PCM である。また、第1及び第2フレーム化部 11, 12 でフレーム化された信号のフレーム長は、例えば 512 サンプルフレーム長である。

40

【0082】

例えば、音声信号  $x_2(t)$  を第2フレーム化部 12 でフレーム化してから音入力オン/オフ制御部 15 に出力することで、結果的に、音声信号処理装置 10 から出力される音声信号  $x_2(t)$  もフレーム化されているものとなり、これにより、音声信号処理装置 10 から出力される音声信号  $x_2(t)$  を利用する音声アプリケーションでは、解りやすいフレーム化された音声信号  $x_2(t)$  で処理をすることができるようになる。

50

## 【0083】

また、前述の実施形態では、検出対象音が人間が発する発話音である場合を説明したが、検出対象音は、人間以外の物体が発する音でもよい。

また、前述の実施形態の説明において、相互相関関数計算部13又はパワースペクトラム計算部21、22及びパワー比計算部23は、無指向性マイクに入力された音信号と、単一指向性マイクに入力された音信号とを比較する比較手段を実現しており、音声/非音声判定部14、24、31は、比較手段の比較結果に基づいて、検出対象音を検出する検出対象音検出手段又は発話音の発話区間を検出する発話区間検出手段を実現している。

## 【0084】

また、前述の実施形態の音声信号処理装置10を音声認識装置に適用することができる。この場合、音声認識装置は、前述したような音声信号処理装置10の構成に加えて、音声信号処理装置10が検出した発話区間の音声信号について音声認識処理をする音声認識処理手段を備える。 10

ここで、音声認識技術としては、例えば、旭化成株式会社が提供する音声認識技術「VORERO」（商標）（<http://www.asahi-kasei.co.jp/vorero/jp/vorero/feature.html>参照）等があり、このような音声認識技術の用いた音声認識装置に適用することもできる。

## 【0085】

また、前述の実施形態の音声信号処理装置10をコンピュータで実現することができる。そして、前述したような音声信号処理装置10の処理内容をコンピュータが所定のプログラムにより実現する。この場合、プログラムは、無指向性マイクで受音した発話音及び雑音の音声信号と単一指向性マイクで受音した前記発話音又は前記雑音のいずれか一方の音声信号とを比較し、その比較結果に基づいて、前記発話音の発話区間を検出する処理をコンピュータに実行させるプログラムになる。 20

## 【図面の簡単な説明】

## 【0086】

【図1】本発明の実施形態の音声信号処理装置を含むシステム全体の構成を示すブロック図である。

【図2】本発明の第1の実施形態におけるマイクの配置を示す図である。

【図3】本発明の第1の実施形態の音声信号処理装置の構成を示すブロック図である。 30

【図4】本発明の第2の実施形態におけるマイクの配置を示す図である。

【図5】本発明の第3の実施形態の音声信号処理装置の構成を示すブロック図である。

【図6】本発明の第5の実施形態の音声信号処理装置の構成を示すブロック図である。

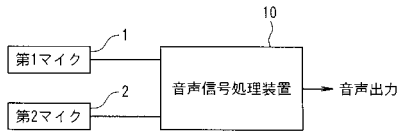
【図7】前記第1の実施形態の他の構成例を示すブロック図である。

## 【符号の説明】

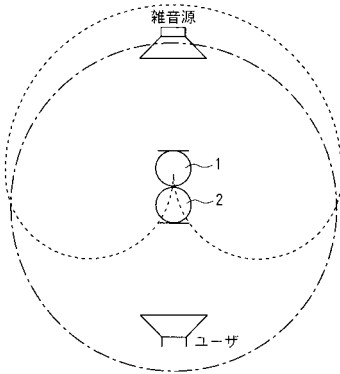
## 【0087】

- 1, 2   マイク
- 10    音声信号処理装置
- 11, 12   フレーム化部
- 13    相互相関関数計算部
- 14, 24, 31   音声/非音声判定部
- 15    音入力オン/オフ制御部
- 21, 22   パワースペクトラム計算部
- 23    パワー比計算部

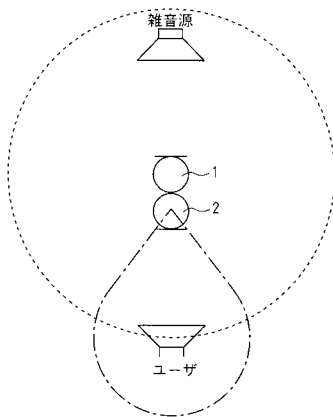
【 図 1 】



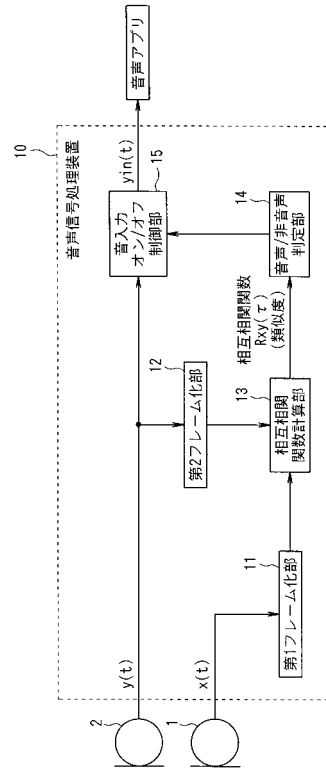
【 図 2 】



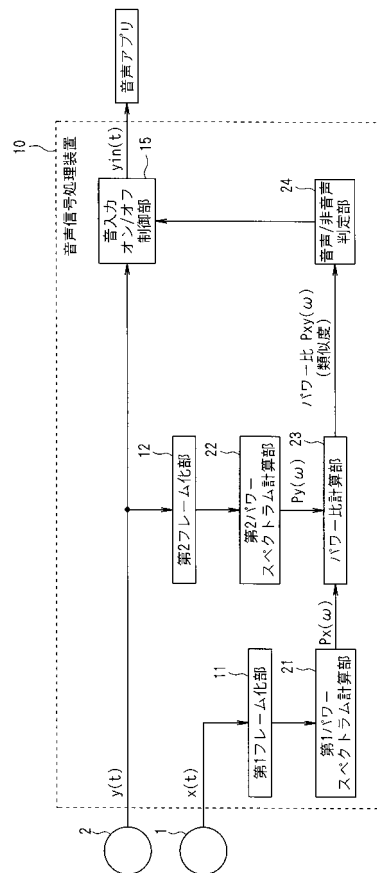
【 図 4 】



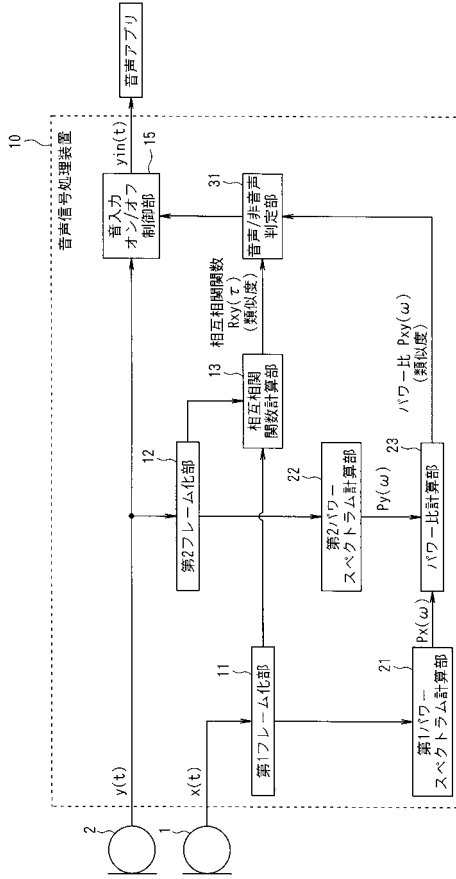
【 図 3 】



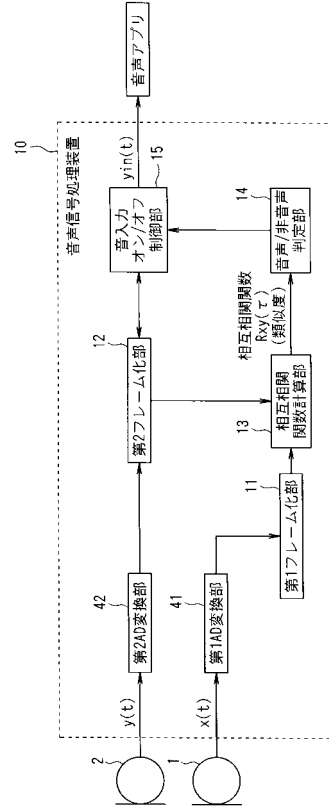
【 図 5 】



【図6】



【図7】



---

フロントページの続き

(51)Int.Cl.<sup>7</sup>

G 1 0 L 21/02

F I

テーマコード(参考)