

(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2005-227512

(P2005-227512A)

(43) 公開日 平成17年8月25日(2005.8.25)

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

G10L 15/20  
G10L 11/02  
G10L 15/00  
G10L 15/02  
G10L 15/04

F I

G10L 3/02 301Z  
H04R 3/00 320  
G10L 3/00 511  
G10L 3/02 301E  
G10L 3/00 513B

テーマコード(参考)

5D015  
5D020

審査請求 未請求 請求項の数 12 O L (全 26 頁) 最終頁に続く

(21) 出願番号 特願2004-35619(P2004-35619)

(22) 出願日 平成16年2月12日(2004.2.12)

(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

(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

Fターム(参考) 5D015 AA06 BB02 CC17 DD02 EE05

FF07 KK02

5D020 BB07

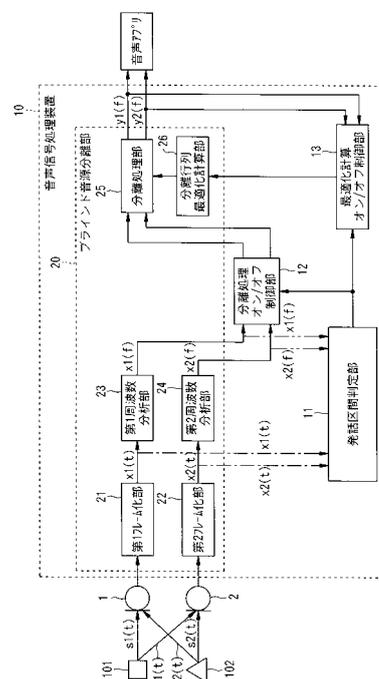
(54) 【発明の名称】 音声信号処理方法及びその装置、音声認識装置並びにプログラム

(57) 【要約】

【課題】 ブラインド音源分離をリアルタイムで行うことを可能にする。

【解決手段】 音声信号処理装置10は、話者音源101からの音声と雑音源102からの音との混合音が入力される第1及び第2マイク1,2と、話者音源101から出力された発話区間を検出する発話区間判定部11と、発話区間判定部11が発話区間を検出した場合、第1及び第2マイク1,2に入力された音声信号 $x_1(t)$ 、 $x_2(t)$ を用いて、分離行列を最適化する分離処理オン/オフ制御部12、最適化計算オン/オフ制御部13及び分離行列最適化計算部26と、分離行列最適化計算部26が最適化した分離行列を用いて、混合音から話者音源101からの音と雑音源102からの音とを分離する分離処理部25とを備える。

【選択図】 図2



**【特許請求の範囲】****【請求項 1】**

検出対象音源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該最適化した分離行列を用いて、前記混合音から検出対象音源からの音と雑音源からの音とを分離するブラインド音源分離 ( B B S :BlindSource Separation ) を行う音信号処理方法において、

前記検出対象音源からの音の検出の有無により前記分離行列の最適化の実行を切替可能にするとともに、前記検出対象音源からの音を検出した場合、前記分離行列の最適化を行い、前記検出対象音源からの音を検出できない場合、前記分離行列の最適化を行わないことを特徴とする音信号処理方法。

10

**【請求項 2】**

前記検出対象音源からの音が所定長以上の音の場合、前記分離行列の最適化を行い、前記検出対象音源からの音が所定長未満の音の場合、前記分離行列の最適化を行わないことを特徴とする請求項 1 記載の音信号処理方法。

**【請求項 3】**

前記ブラインド音源分離では、無指向性マイクに前記混合音が入力され、単一指向性マイクに前記検出対象音源からの音又は前記雑音源からの音のいずれか一方が入力され、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音の音信号を用いて前記分離行列の最適化を行い、かつ当該最適化した分離行列を用いて前記混合音から検出対象音源からの音と雑音源からの音とを分離しており、

20

前記無指向性マイクに入力された混合音の音信号と前記単一指向性マイクに入力された音の音信号とを比較し、その比較結果に基づいて、前記検出対象音源からの音を検出することを特徴とする請求項 1 又は 2 記載の音信号処理方法。

**【請求項 4】**

検出対象音源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該最適化した分離行列を用いて、前記混合音から検出対象音源からの音と雑音源からの音とを分離するブラインド音源分離 ( B B S :BlindSource Separation ) を行う音信号処理装置において、

前記混合音が入力される第 1 マイクと、

前記検出対象音源からの音と雑音源からの音とのうちの少なくとも一方が入力される第 2 マイクと、

30

前記検出対象音源からの音を検出する対象音検出手段と、

前記対象音検出手段が検出対象音を検出した場合、前記第 1 及び第 2 マイクに入力された音の音信号を用いて前記分離行列を最適化する分離行列最適化手段と、

前記分離行列最適化手段が最適化した分離行列を用いて、前記第 1 マイクに入力された混合音から検出対象音源からの音と雑音源からの音とを分離する分離手段と、

を備えることを特徴とする音信号処理装置。

**【請求項 5】**

前記分離行列最適化手段は、前記検出対象音検出手段が検出した検出対象音が所定長以上の音の場合、前記分離行列の最適化を行うことを特徴とする請求項 4 記載の音信号処理装置。

40

**【請求項 6】**

前記第 1 マイクは、前記混合音が入力されるように配置された無指向性マイクであり、第 2 のマイクは、前記検出対象音源からの音と雑音源からの音とのうちのいずれか一方が入力されるように配置された単一指向性マイクであり、

前記対象音検出手段は、前記第 1 マイクに入力された混合音の音信号と前記第 2 マイクに入力された音の音信号とを比較して、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 4 又は 5 記載の音信号処理装置。

**【請求項 7】**

前記第 1 マイクに入力された混合音の音信号及び第 2 マイクに入力された検出対象音源

50

からの音と雑音源からの音とのうちのいずれか一方の音の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、

前記対象音検出手段は、前記フレーム化手段から出力されるフレーム単位で、前記第1マイクに入力された混合音の音信号と、前記第2マイクに入力された音の音信号とを比較し、そ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6記載の音声認識装置。

【請求項8】

発話源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該最適化した分離行列を用いて、前記混合音から発話源からの音と雑音源からの音とを、ブラインド音源分離(BBS:BlindSource Separation)により分離し、その分離した発話源

10

からの音について音声認識処理を行う音声認識装置において、

前記混合音が入力される第1マイクと、  
前記発話源からの音と雑音源からの音とのうちの少なくとも一方が入力される第2マイクと、

前記発話源からの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発話区間を検出した場合、前記第1及び第2マイクに入力された音信号を用いて前記分離行列を最適化する分離行列最適化手段と、

前記分離行列最適化手段が最適化した分離行列を用いて、前記第1マイクに入力された混合音から発話源からの音と雑音源からの音とを分離する分離手段と、

前記分離手段が分離した発話源からの音について、音声認識処理を行う音声認識処理手

20

段と、  
を備えることを特徴とする音声認識装置。

【請求項9】

前記分離行列最適化手段は、前記発話区間検出手段が発話区間が所定長以上の場合、前記分離行列の最適化を行うことを特徴とする請求項8記載の音声認識装置。

【請求項10】

前記第1マイクは、前記混合音が入力されるように配置された無指向性マイクであり、第2のマイクは、前記発話源からの音と雑音源からの音とのうちのいずれか一方が入力されるように配置された単一指向性マイクであり、

前記発話区間検出手段は、前記第1マイクに入力された混合音の音信号と前記第2マイクに入力された音の音信号とを比較して、その比較結果に基づいて、前記発話区間を検出することを特徴とする請求項8又は9記載の音声認識装置。

30

【請求項11】

前記第1に入力された混合音の音信号及び第2マイクに入力された発話源からの音と雑音源からの音とのうちのいずれか一方の音の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、

前記発話区間検出手段は、前記フレーム化手段から出力されるフレーム単位で、前記第1マイクに入力された混合音の音信号と、前記第2マイクに入力された音の音信号とを比較し、そ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。

40

【請求項12】

検出対象音源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該分離行列を用いて、前記混合音から検出対象音源からの音と雑音源からの音とを分離するブラインド音源分離(BBS:BlindSource Separation)をコンピュータに実現させるプログラムにおいて、

前記検出対象音源からの音の検出の有無により前記分離行列の最適化の実行を切替可能にするとともに、前記検出対象音源からの音を検出した場合、前記分離行列の最適化を行い、前記検出対象音源からの音を検出できない場合、前記分離行列の最適化を行わないように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】

50

## 【技術分野】

## 【0001】

本発明は、音信号処理方法、音信号処理装置、音声認識装置及びプログラムに関し、特に混合音から検出対象音を分離して取り出すブラインド音源分離（BBS:BlindSource Separation）が適合される音信号処理方法、音信号処理装置、音声認識装置及びプログラムに関する。

## 【背景技術】

## 【0002】

ブラインド音源分離（BBS:Blind Source Separation）では、複数チャンネルに入力された混合音を用いて、独立成分分析（ICA:IndependentComponent Analysis）の技術により、分離行列を最適化（学習）する。これにより、分離行列が目的とする音を分離する最適解に近づく。そして、ブラインド音源分離では、そのように最適化した分離行列を用いて、混合音から目的の音を分離して取り出している。ここで、混合音として、話者音源（発話源）からの音（発話）と雑音源からの音が混ざり合った音が挙げられ、このような場合、分離目的の音は、話者音源からの音（発話）になる。

## 【発明の開示】

## 【発明が解決しようとする課題】

## 【0003】

リアルタイムでブラインド音源分離をする場合、混合音で分離行列を最適化しつつ、混合音から目的の音を分離するような態様となる。このようにリアルタイムでブラインド音源分離を実現する場合には、混合音中に分離目的の音が断続的又は不規則に含まれるようになる。

しかし、従来のブラインド音源分離のシステムは、オフライン処理によりブラインド音源分離をすることを前提としている。すなわち、従来のブラインド音源分離のシステムは、分離目的の音が連続して入力される場合を前提とし、その前提の下、分離行列を最適化しつつ、その最適化した分離行列で目的の音を分離するように構成されている。このことから、従来のシステムでリアルタイムでブラインド音源分離を行うと、混合音に分離目的の音が断続的又は不規則に含まれる結果、分離目的の音の特定が困難になることから、分離行列を最適化できなくなる。この結果、目的の音を高精度で分離できなくなる。

## 【0004】

このように、従来のシステムは、オフライン処理でブラインド音源分離を行う必要があった。このような結果、従来のシステムは、実用性に欠けたものとなっていた。

さらに、従来のシステムでは、分離行列の最適化処理中に分離目的としない他の音が長時間継続して入力されてしまうと、当該他の音で分離行列を最適化してしまう。この場合、分離行列が間違った局所最適解に落ち込んでしまう。このように分離行列が間違った局所最適解に落ち込んでしまうと、その後、分離目的の音を入力しても、分離行列が最適化しなくなってしまう。

本発明は、前記問題に鑑みてなされたものであり、ブラインド音源分離をリアルタイムで行うことを可能にする音信号処理方法、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。

## 【課題を解決するための手段】

## 【0005】

請求項1記載の音信号処理方法は、検出対象音源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該最適化した分離行列を用いて、前記混合音から検出対象音源からの音と雑音源からの音とを分離するブラインド音源分離（BBS:BlindSource Separation）を行う音信号処理方法である。この音信号処理方法は、前記検出対象音源からの音の検出の有無により前記分離行列の最適化の実行を切替可能にするとともに、前記検出対象音源からの音を検出した場合、前記分離行列の最適化を行い、前記検出対象音源からの音を検出できない場合、前記分離行列の最適化を行わないことを特徴とする。なお、検出対象音源からの音には、人間が発する発話音の他、物体が発する音も含まれ

る。

【0006】

また、請求項2記載の音信号処理方法は、請求項1記載の音信号処理方法において、前記検出対象音源からの音が所定長以上の音の場合、前記分離行列の最適化を行い、前記検出対象音源からの音が所定長未満の音の場合、前記分離行列の最適化を行わないことを特徴とする。

また、請求項3記載の音信号処理方法は、請求項1又は2記載の音信号処理方法において、前記ブラインド音源分離では、無指向性マイクに前記混合音が入力され、単一指向性マイクに前記検出対象音源からの音又は前記雑音源からの音のいずれか一方が入力され、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音の音信号を用いて前記分離行列の最適化を行い、かつ当該最適化した分離行列を用いて前記混合音から検出対象音源からの音と雑音源からの音とを分離しており、前記無指向性マイクに入力された混合音の音信号と前記単一指向性マイクに入力された音の音信号とを比較し、その比較結果に基づいて、前記検出対象音源からの音を検出することを特徴とする。

10

20

30

40

50

【0007】

また、請求項4記載の音信号処理装置は、検出対象音源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該最適化した分離行列を用いて、前記混合音から検出対象音源からの音と雑音源からの音とを分離するブラインド音源分離 (BBS: Blind Source Separation) を行う音信号処理装置である。この音信号処理装置は、前記混合音が入力される第1マイクと、前記検出対象音源からの音と雑音源からの音とのうちの少なくとも一方が入力される第2マイクと、前記検出対象音源からの音を検出する対象音検出手段と、前記対象音検出手段が検出対象音を検出した場合、前記第1及び第2マイクに入力された音の音信号を用いて前記分離行列を最適化する分離行列最適化手段と、前記分離行列最適化手段が最適化した分離行列を用いて、前記第1マイクに入力された混合音から検出対象音源からの音と雑音源からの音とを分離する分離手段と、を備えることを特徴とする。

【0008】

また、請求項5記載の音信号処理装置は、請求項4記載の音信号処理装置において、前記分離行列最適化手段が、前記検出対象音検出手段が検出した検出対象音が所定長以上の音の場合、前記分離行列の最適化を行うことを特徴とする。

また、請求項6記載の音信号処理装置は、請求項4又は5記載の音信号処理装置において、前記第1マイクが、前記混合音が入力されるように配置された無指向性マイクであり、第2のマイクが、前記検出対象音源からの音と雑音源からの音とのうちのいずれか一方が入力されるように配置された単一指向性マイクであり、前記対象音検出手段が、前記第1マイクに入力された混合音の音信号と前記第2マイクに入力された音の音信号とを比較して、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。

【0009】

また、請求項7記載の音信号処理装置は、請求項6記載の音信号処理装置において、前記第1マイクに入力された混合音の音信号及び第2マイクに入力された検出対象音源からの音と雑音源からの音とのうちのいずれか一方の音の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、前記対象音検出手段が、前記フレーム化手段から出力されるフレーム単位で、前記第1マイクに入力された混合音の音信号と、前記第2マイクに入力された音の音信号とを比較し、そ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。

【0010】

また、請求項8記載の音声認識装置は、発話源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該最適化した分離行列を用いて、前記混合音から発話源からの音と雑音源からの音とを、ブラインド音源分離 (BBS: Blind Source Separation) により分離し、その分離した発話源からの音について音声認識処理を行う音声認識装置である。この音声認識装置は、前記混合音が入力される第1マイクと、前記発話源から

の音と雑音源からの音とのうちの少なくとも一方が入力される第2マイクと、前記発話源からの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発話区間を検出した場合、前記第1及び第2マイクに入力された音信号を用いて前記分離行列を最適化する分離行列最適化手段と、前記分離行列最適化手段が最適化した分離行列を用いて、前記第1マイクに入力された混合音から発話源からの音と雑音源からの音とを分離する分離手段と、前記分離手段が分離した発話源からの音について、音声認識処理を行う音声認識処理手段と、を備える。

**【0011】**

また、請求項9記載の音声認識装置は、請求項8記載の音声認識装置において、前記分離行列最適化手段が、前記発話区間検出手段が検出した発話区間が所定長以上の場合、前記分離行列の最適化を行うことを特徴とする。

10

また、請求項10記載の音声認識装置は、請求項8又は9記載の音声認識装置において、前記第1マイクが、前記混合音が入力されるように配置された無指向性マイクであり、第2のマイクが、前記発話源からの音と雑音源からの音とのうちのいずれか一方が入力されるように配置された単一指向性マイクであり、前記発話区間検出手段が、前記第1マイクに入力された混合音の音信号と前記第2マイクに入力された音の音信号とを比較して、その比較結果に基づいて、前記発話区間を検出することを特徴とする。

**【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、前記第1に入力された混合音の音信号及び第2マイクに入力された発話源からの音と雑音源からの音とのうちのいずれか一方の音の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、前記発話区間検出手段が、前記フレーム化手段から出力されるフレーム単位で、前記第1マイクに入力された混合音の音信号と、前記第2マイクに入力された音の音信号とを比較し、そ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記載のプログラムは、検出対象音源からの音と雑音源からの音との混合音により分離行列を最適化していき、当該分離行列を用いて、前記混合音から検出対象音源からの音と雑音源からの音とを分離するブラインド音源分離(BBS:BlindSource Separation)をコンピュータに実現させるプログラムである。このプログラムは、前記検出対象音源からの音の検出の有無により前記分離行列の最適化の実行を切替可能にするとともに、前記検出対象音源からの音を検出した場合、前記分離行列の最適化を行い、前記検出対象音源からの音を検出できない場合、前記分離行列の最適化を行わないようにコンピュータに実行させることを特徴とする。

30

**【発明の効果】****【0014】**

本発明によれば、発話源からの音を検出した場合、分離行列の最適化を行い、前記発話源からの音を検出できない場合、分離行列の最適化を行わないので、断続的又は不規則にシステムに入力される発話源からの音に対してのみ分離行列の最適化を行うことができる。これにより、リアルタイムでブラインド音源分離を行うことができる。

40

また、請求項2、5及び9記載の発明によれば、検出対象音源からの音又は発話源からの音が所定長以上の場合、分離行列の最適化を行うようにすることで、検出対象音源からの音又は発話源からの音に対して最適解の分離行列を得ることができる。

**【0015】**

また、請求項3、6及び10記載の発明によれば、無指向性マイクで検出対象音源からの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源からの音(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、前記検出対象音源からの音(発話源からの音)を検出することができる。これにより、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の構築が可能になる。

50

## 【発明を実施するための最良の形態】

## 【0016】

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。

図2は音声信号処理装置10の構成を示す。

図2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21, 22、第1及び第2周波数分析部23, 24、分離処理部25、分離行列最適化計算部26、発話区間判定部11、分離処理オン/オフ制御部12及び最適化計算オン/オフ制御部13を備えている。 10

## 【0017】

なお、第1及び第2フレーム化部21, 22、第1及び第2周波数分析部23, 24、分離処理部25及び分離行列最適化計算部26は、ブラインド音源分離（BBS:Blind Source Separation）を実現するブラインド音源分離部20を構成している。すなわち、ブラインド音源分離部20は、このような構成を備えることで、複数チャンネルに入力された混合音により、独立成分分析（ICA:Independent Component Analysis）の技術を用いて分離行列を最適化する一方、当該最適化した分離行列を用いて、前記混合音から話者音源（発話源）からの音（発話音）と雑音源からの音（雑音）とを分離する音声信号処理を実現する。 20

## 【0018】

このような音声信号処理装置10の構成において、第1及び第2マイク1, 2から入力された2chの音声信号 $x_1(t)$ ,  $x_2(t)$ はそれぞれ、第1及び第2フレーム化部21, 22に入力される。

ここで、音声信号 $x_1(t)$ ,  $x_2(t)$ は、話者音源（発話源）101が発した音 $s_1(t)$ と雑音源102が発した音 $s_2(t)$ とが混ざり合った混合音信号である。雑音 $s_2(t)$ としては、話者音源の周囲の音、話者音源以外の他の者の音声等が挙げられる。

## 【0019】

第1フレーム化部21では、第1マイク1から入力された音声信号 $x_1(t)$ 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 $x_1(t)$ を第1周波数分析部23に出力する。第2フレーム化部22では、第2マイク2から入力される音声信号 $x_2(t)$ 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 $x_2(t)$ を第2周波数分析部24に出力する。ここでは、第1及び第2フレーム化部21, 22は、入力されてくる音声信号 $x_1(t)$ ,  $x_2(t)$ 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。 30

## 【0020】

第1及び第2周波数分析部23, 24はそれぞれ、フレーム単位で音声信号 $x_1(t)$ ,  $x_2(t)$ をFFT（Fast Fourier Transform）により周波数分析して、観測信号（Observed signals） $x_1(f)$ ,  $x_2(f)$ を生成し、その観測信号 $x_1(f)$ ,  $x_2(f)$ を分離処理オン/オフ制御部12に出力する。 40

なお、観測信号 $x_1(f)$ ,  $x_2(f)$ とは、当該ブラインド音源分離（BBS:Blind Source Separation）の技術において、混合音の分離を行う分離行列に入力される信号のことをいう。

## 【0021】

分離処理オン/オフ制御部12は、発話区間判定部11からの発話区間判定結果（制御信号）に基づいて、第1及び第2周波数分析部23, 24それぞれからの観測信号 $x_1(f)$ ,  $x_2(f)$ を後段の分離処理部25に出力する。

発話区間判定部11は、第1及び第2マイク1, 2から入力された音声信号 $x_1(t)$  50

、 $x_2(t)$ に基づいて、当該音声信号 $x_1(t)$ 、 $x_2(t)$ に含まれている発話音声の区間（発話区間）を判定するように構成されている。例えば、発話区間判定部11は、第1及び第2マイク1, 2から入力された音声信号 $x_1(t)$ 、 $x_2(t)$ 、具体的には第1及び第2フレーム化部11, 12から出力されたフレーム単位の音声信号 $x_1(t)$ 、 $x_2(t)$ 又は第1及び第2周波数分析部23, 24から出力された信号 $x_1(f)$ 、 $x_2(f)$ に基づいて、当該フレーム単位で発話区間の判定を行う。具体的には、発話区間判定部11は、所定長（所定時間）以上の発話区間を検出したときに、発話区間を検出した旨の信号を判定結果（制御信号）として、分離処理オン/オフ制御部12及び最適化計算オン/オフ制御部13に出力する。なお、発話区間判定部11の具体的な構造については、後述する第2乃至第4の実施形態として説明する。

10

#### 【0022】

これにより、分離処理オン/オフ制御部12は、発話区間判定部11から発話区間を検出した結果が入力された場合、分離処理部25のオン制御として、第1及び第2周波数分析部23, 24それぞれからの観測信号 $x_1(f)$ 、 $x_2(f)$ を分離処理部25に出力する。また、分離処理オン/オフ制御部12は、発話区間判定部11が発話区間を検出していない場合、分離処理部25のオフ制御として、第1及び第2周波数分析部23, 24それぞれからの観測信号 $x_1(f)$ 、 $x_2(f)$ を分離処理部25に出力しない。このとき、分離処理オン/オフ制御部12から分離処理部25への観測信号 $x_1(f)$ 、 $x_2(f)$ の出力のオン及びオフは、発話区間判定部11が発話区間を検出したフレームに対応するフレームを単位として行う。

20

#### 【0023】

分離処理部25は、分離行列最適化計算部26により最適化された分離行列により、観測信号 $x_1(f)$ 、 $x_2(f)$ から分離信号 $y_1(f)$ 、 $y_2(f)$ を分離抽出する。そして、分離処理部25は、音声信号 $s_1(t)$ 、 $s_2(t)$ とされる分離信号 $y_1(f)$ 、 $y_2(f)$ を後段に出力する。

分離行列最適化計算部26は、分離処理部25が得た分離信号 $y_1(f)$ 、 $y_2(f)$ が入力されており、この分離信号 $y_1(f)$ 、 $y_2(f)$ に基づく分離行列の最適化処理として、最適解の分離行列を得る。そして、分離行列最適化計算部26は、その最適化した分離行列を分離処理部25に出力する。すなわち、分離処理部25は、当該分離処理部25が得る分離信号 $y_1(f)$ 、 $y_2(f)$ を用いて分離行列最適化計算部26で最適化された分離行列を用いて、それ以降に当該分離処理部25に入力される観測信号 $x_1(f)$ 、 $x_2(f)$ から分離信号 $y_1(f)$ 、 $y_2(f)$ を分離抽出しているのである。

30

#### 【0024】

一方、分離行列最適化計算部26は、最適化計算オン/オフ制御部13によりオン及びオフ制御がなされる。具体的には、最適化計算オン/オフ制御部13は、発話区間判定部11から発話区間を検出した結果が入力された場合、分離行列最適化計算部26をオン制御しており、分離行列最適化計算部26はこのオン制御により、分離処理部25が出力した分離信号 $y_1(f)$ 、 $y_2(f)$ に基づいて、分離行列の最適化処理を実施する。また、最適化計算オン/オフ制御部13は、発話区間判定部11が発話区間を検出していない場合、分離行列最適化計算部26をオフ制御しており、分離行列最適化計算部26はこの

40

#### 【0025】

以上のように音声信号処理装置10が構成されている。

次に図3を用いて、第1及び第2マイク1, 2から入力された2chの音声信号（混合音声信号） $x_1(t)$ 、 $x_2(t)$ に対する処理に沿って、音声信号処理装置10の一連の動作を説明する。なお、ここでの動作は、分離行列を最適化（学習）する際の動作になる。

第1及び第2マイク1, 2からの音声信号 $x_1(t)$ 、 $x_2(t)$ は、第1及び第2フレーム化部21, 22に入力される。

第1及び第2フレーム化部21, 22は、各音声信号 $x_1(t)$ 、 $x_2(t)$ をフレ

50

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  $x_1(t)$  ,  $x_2(t)$  を第 1 及び第 2 周波数分析部 23 , 24 に出力する（ステップ S1）。

【0026】

第 1 及び第 2 周波数分析部 23 , 24 では、フレーム単位で、音声信号  $x_1(t)$  ,  $x_2(t)$  から観測信号  $x_1(f)$  ,  $x_2(f)$  を生成し、その観測信号  $x_1(f)$  ,  $x_2(f)$  を分離処理オン/オフ制御部 12 に出力する（ステップ S2）。

一方、発話区間判定部 11 は、第 1 及び第 2 マイク 1 , 2 から入力された音声信号  $x_1(t)$  ,  $x_2(t)$  中の発話区間の判定をフレーム単位で行い（ステップ S3）、発話区間（発話フレーム）を検出する（ステップ S4）。そして、発話区間判定部 11 は、発話区間を検出した場合、当該発話区間が最短発話長以上か否かを判定する（ステップ S5）  
10  
ここで、発話区間判定部 11 は、発話区間が最短発話長以上の場合、発話区間を検出した旨の判定結果を分離処理オン/オフ制御部 12 及び最適化計算オン/オフ制御部 13 に出力する。また、発話区間判定部 11 は、発話区間を検出できなかった場合、又は発話区間は検出できたが、その発話区間が最短発話長未満であった場合、発話区間を検出できなかったとして、その旨の判定結果を分離処理オン/オフ制御部 12 及び最適化計算オン/オフ制御部 13 に出力する。

【0027】

分離行列最適化計算部 26 は、分離処理部 25 から分離行列を読み出す（ステップ S6）。そして、分離行列最適化計算部 26 は、その読み出した分離行列の最適化計算を行う（ステップ S7）。具体的には次のような処理により分離行列の最適化計算を行う。  
20

分離処理オン/オフ制御部 12 では、発話区間判定部 11 が発話区間を検出した場合、第 1 及び第 2 周波数分析部 23 , 24 それぞれからの観測信号  $x_1(f)$  ,  $x_2(f)$  を後段の分離処理部 25 に出力する。そして、分離処理部 25 は、最新の分離行列により観測信号  $x_1(f)$  ,  $x_2(f)$  から分離信号  $y_1(f)$  ,  $y_2(f)$  を得る。

【0028】

その一方で、最適化計算オン/オフ制御部 13 は、発話区間判定部 11 が発話区間を検出した場合、分離行列最適化計算部 26 をオン制御する。分離行列最適化計算部 26 は、オン制御により、分離処理部 25 が得た分離信号  $y_1(f)$  ,  $y_2(f)$  を取り込み、この分離信号  $y_1(f)$  ,  $y_2(f)$  に基づいて前記読み出した分離行列を最適化する。

このように分離行列最適化計算部 26 で分離行列の最適化計算を行う。そして、分離行列最適化計算部 26 は、その最適化した分離行列を分離処理部 25 に出力し、分離処理部 25 は、その分離行列を保存する（ステップ S8）。  
30

【0029】

そして、分離処理部 25 は、このように最適化された最新の分離行列を用いて、観測信号  $x_1(f)$  ,  $x_2(f)$  から分離信号  $y_1(f)$  ,  $y_2(f)$  を得る（ステップ S9）。

このように音声信号処理装置 10 は、分離処理部 25 で得た分離信号  $y_1(f)$  ,  $y_2(f)$  を例えば音声アプリケーションに出力する。

音声アプリケーションは、例えば音声を認識して各種処理を行うアプリケーションである。例えば、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。このような音声アプリケーションは、話者音源（発話源）101 が発した音声信号  $s_1(t)$  である分離信号  $y_1(f)$  に基づいて、音声を認識して、所定の処理を行う。  
40

【0030】

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、音声信号処理装置 10 は、発話区間を検出した場合にのみ、分離行列の最適化計算を行っている。これにより、分離目的の音である発話音源からの音が音声信号処理装置 10 に断続的又は不規則に輸入されてくる場合でも、音声信号処理装置 10 は、分離行列を最適化することができる。これにより、分離目的の音である発話音源からの音が音声信号処理装置 10 に断続的又は不規則に輸入されてくる場合でも、音声信号処理  
50

装置 10 は、目的の音である発話音源からの音を高精度で分離できるようになる。このように、音声信号処理装置 10 は、リアルタイムでブラインド音源分離を実現できるようになり、実用性に優れたものとなる。

#### 【0031】

また、このように発話区間を検出した場合にのみ分離行列の最適化計算を行うようにすることで、分離目的外の音が入力されても分離行列の最適化計算が行われないので、そのような分離目的外の音により分離行列が間違った局所最適解に落ち込んでしまうようなことを防止できる。

また、前述したように、音声信号処理装置 10 は、発話区間が最短発話長以上の場合に限って、分離行列の最適化計算を行っている。一般的には、ブラインド音源分離のシステムに入力される音（学習対象の音）がある一定以上の長さがあると、分離行列の最適化は良好となる。このようなことから、音声信号処理装置 10 は、発話区間が最短発話長以上の場合に限って分離行列の最適化計算を行うようにすることで、分離目的の音に最適解の分離行列を得ることができるようになる。なお、音声信号処理装置 10 が分離した音声を音声認識システム（音声アプリケーション）が利用するとした場合、前記一定以上の長さとは、例えばコマンド最短長さや、1 発話最短長さとなる。

#### 【0032】

そして、このように音声信号処理装置 10 では、高精度で目的の音声を分離できるので、このように音声信号処理装置 10 が分離した音声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。

#### 【0033】

また、前述したように、発話区間判定部 11 が発話区間を検出した場合には、分離処理オン/オフ制御部 12 が第 1 及び第 2 周波数分析部 23, 24 それぞれからの観測信号  $x_1(f)$ ,  $x_2(f)$  を後段の分離処理部 25 に出力する一方で、最適化計算オン/オフ制御部 13 が分離行列最適化計算部 26 をオン制御して、分離行列最適化計算部 26 に分離行列の最適化処理を実施させている。

#### 【0034】

よって、発話区間判定部 11 が発話区間を検出した場合にのみ、分離行列の最適化計算を行うのであれば、分離処理オン/オフ制御部 12 と最適化計算オン/オフ制御部 13 とのいずれか一方を備えるだけでよいといえる。しかし、分離処理オン/オフ制御部 12 や最適化計算オン/オフ制御部 13 の応答性を考慮して、これら両方をシステムに備えることで、それら構成要素の特性に対するロバスト性を上げて、分離行列の最適化処理を行うことができるようになる。

#### 【0035】

また、前述したように、音声信号  $x_1(t)$ ,  $x_2(t)$  を第 1 及び第 2 フレーム化部 21, 22 でフレーム化したものを、分離処理オン/オフ制御部 12 及び分離処理部 25 に出力している。このようにすることで、結果的に、音声信号処理装置 10 から出力される分離信号  $y_1(f)$  である音声信号  $s_1(t)$  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  $s_1(t)$  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  $s_1(t)$  で処理をすることができるようになる。

#### 【0036】

ここで、図 4 を用いて効果を説明する。

図 4 中 (A) は、オフラインによりブラインド音源分離を行う場合を示し（従来の手法）、図 4 中 (B) 及び (C) は、リアルタイムでブラインド音源分離を行う場合を示す。

従来の手法をそのまま適用して、リアルタイムでブラインド音源分離をしてしまうと、図 4 中 (B) に従来法として示すように、システムに雑音のみが入力されている場合でも、その雑音により分離行列を最適化してしまう。この場合、分離行列が劣化してしまう。

10

20

30

40

50

この結果、最適化された分離行列では、目的とする信号（音声信号）を分離することができなくなる（結果不明となる）。

【0037】

一方、本発明を適用して、リアルタイムでブラインド音源分離をした場合、図4中（B）に本発明法として示すように、システムに雑音のみが入力されているときには、分離行列の最適化は実施されず、システムに雑音と目的とする信号（音声信号）とが入力されたときに、分離行列の最適化は実施される。この結果、最適化された分離行列により、雑音とともに入力されてきた目的とする信号（音声信号）を精度よく分離することができる。

【0038】

また、従来の手法をそのまま適用して、リアルタイムでブラインド音源分離をした場合、システムに雑音のみ又は雑音と分離目的外の信号とが混じり合い、長時間入力されると、図4中（C）に従来法として示すように、分離行列が間違っただけの局所最適解に落ち込んでしまう。この結果、目的とする信号（音声信号）を分離することができなくなる（結果不明となる）。

しかし、本発明を適用した場合には、図4中（C）に本発明法として示すように、システムに雑音と目的とする信号（音声信号）とが入力されたときに分離行列の最適化を実施するので、そのように分離行列が間違っただけの局所最適解に落ち込んでしまうことを防止できる。

【0039】

次に第2の実施形態を説明する。

この第2の実施形態は、発話区間判定部11を具体的な構成とした音声信号処理装置10であり、発話区間判定部11が、第1及び第2マイク1, 2で受音した音声信号 $x_1(t)$ ,  $x_2(t)$ の相関度により発話区間を検出するように構成されている。

図5は、その第2の実施形態における発話区間判定部11の構成を示し、図6は、発話区間判定部11の構成に対応する第1及び第2マイクの配置を示す。

【0040】

この第2の実施形態では、第1マイク1として単一指向性マイクを使用し、第2マイク2として無指向性マイクを使用している。そして、第1及び第2マイク1, 2は、図6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をその指向方向が発話音源（ユーザ）の位置に対して反対側となるように配置する。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6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6に示す一点鎖線は、第2マイク2の指向特性を示す。

【0041】

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 $s_2(t)$ は、第1及び第2マイク1, 2で受音でき、発話音源（ユーザ）からの音 $s_1(t)$ は第2マイク2だけが受音できるようになる。

このように配置した第1及び第2マイク1, 2から入力された音声信号 $x_1(t)$ ,  $x_2(t)$ はそれぞれ、前述したように、第1及び第2フレーム化部21, 22に入力される。そして、前述したように、第1フレーム化部21では、第1マイク1から入力された音声信号 $x_1(t)$ をフレーム化（或いはフレーム分割）し、また、第2フレーム化部22では、第2マイク2から入力される音声信号 $x_2(t)$ をフレーム化（或いはフレーム分割）する。そして、このように各フレーム化部21, 22で複数フレームにされた音声信号 $x_1(t)$ ,  $x_2(t)$ は発話区間判定部11に入力される。

【0042】

発話区間判定部11は、図5に示すように、相互相関関数計算部31及び音声/非音声判定部41を備えている。このような発話区間判定部11において、各フレーム化部21, 22で複数フレームにされた音声信号 $x_1(t)$ ,  $x_2(t)$ が相互相関関数計算部31に入力される。

相互相関関数計算部31は、第1フレーム化部21から出力されるフレームと、第2フ

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  $x_1(t)$  と、第 2 マイク 2 に入力された音声信号  $x_2(t)$  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 31 は、下記 (1) 式により、相互相関関数  $R(\tau)$  を算出する。

【0043】

【数 1】

$$R(\tau) = \frac{1}{T} \int_{-T/2}^{T/2} x_1(t)x_2(t-\tau)dt \quad \dots (1)$$

10

【0044】

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 相当分の遅延時間を与えて、相互相関関数  $R(\tau)$  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  $R(\tau)$  を得ることができ、精度よく相互相関関数  $R(\tau)$  を得ることができる。

20

【0045】

このように算出された相互相関関数  $R(\tau)$  は、相関関係を求める 2 つの音声信号  $x_1(t)$ ,  $x_2(t)$  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  $x_1(t)$ ,  $x_2(t)$  が異なっているほど、0 に近くなる。相互相関関数計算部 31 は、このような相互相関関数  $R(\tau)$  を音声 / 非音声判定部 32 に出力する。

音声 / 非音声判定部 32 は、相互相関関数  $R(\tau)$  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。

【0046】

前述したように、発話音源 (ユーザ) と雑音源に対して図 6 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  $s_2(t)$  を第 1 及び第 2 マイク 1, 2 で受信し、発話音源 (ユーザ) からの音  $s_1(t)$  を第 2 マイク 2 だけで受信している。

一方、相互相関関数  $R(\tau)$  は、前述したように、相関関係を求める 2 つの音声信号  $x_1(t)$ ,  $x_2(t)$  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  $x_1(t)$ ,  $x_2(t)$  が異なっているほど 0 に近くなる。

【0047】

このようなことから、雑音源からの音  $s_2(t)$  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  $R(\tau)$  は大きい値になる。一方、発話音源 (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  $R(\tau)$  は 0 に向かって減少する。

30

40

【0048】

このように、発話音源 (ユーザ) から発話があった場合には相互相関関数  $R(\tau)$  は 0 に向かって減少することから、音声 / 非音声判定部 32 は、相互相関関数  $R(\tau)$  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 32 は、相互相関関数  $R(\tau)$  が判定用しきい値  $r_1$  未満の場合 ( $R(\tau) < r_1$ )、音声区間と判定し、それ以外の場合 ( $R(\tau) \geq r_1$ )、非

50

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 3 2 は、このような判定をフレーム単位で行う。発話区間判定部 1 1 は、このように音声 / 非音声判定部 3 2 で得た音声区間（発話区間）の判定結果を分離処理オン / オフ制御部 1 2 及び最適化計算オン / オフ制御部 1 3 に出力する。

【0049】

以上のように、発話区間判定部 1 1 では、相互相関関数計算部 1 3 が、第 1 及び第 2 フレーム化部 2 1, 2 2 それぞれから出力されるフレーム単位で相互相関関数  $R(\quad)$  を算出して、算出した相互相関関数  $R(\quad)$  を音声 / 非音声判定部 3 2 に出力する。音声 / 非音声判定部 3 2 では、相互相関関数  $R(\quad)$  と判定用しきい値  $r_1$  とを比較し、相互相関関数  $R(\quad)$  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 3 2 は、その判定結果を分離処理オン / オフ制御部 1 2 及び最適化計算オン / オフ制御部 1 3 に出力する。

10

【0050】

そして、分離処理オン / オフ制御部 1 2 は、前述したように、発話区間判定部 1 1 からの発話区間の判定結果に基づいて、分離処理部 2 5 への観測信号  $x_1(f)$ ,  $x_2(f)$  の出力をオン及びオフ制御する。また、最適化計算オン / オフ制御部 1 3 は、前述したように、発話区間判定部 1 1 からの発話区間の判定結果に基づいて、分離行列最適化計算部 2 6 のオン及びオフを制御する。

【0051】

なお、第 1 及び第 2 マイク 1, 2 の配置については、前記図 6 に示した態様に限定されるものではない。例えば、発話音源（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、第 1 及び第 2 マイク 1, 2 を配置してもよい。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 7 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、その指向方向が発話音源（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。なお、図 7 に示す点線は、第 1 マイク 1 の指向特定を示し、図 7 に示す一点鎖線は、発話音源（ユーザ）を基準にした第 2 マイク 2 の指向特性を示す。

20

【0052】

そして、このように第 1 及び第 2 マイク 1, 2 を配置した場合には、相互相関関数計算部 3 1 及び音声 / 非音声判定部 3 2 は次のような計算を行う。

30

発話音源（ユーザ）からの音  $s_1(t)$  を第 1 及び第 2 マイク 1, 2 で受音し、雑音源からの音  $s_2(t)$  を第 1 マイク 1 だけが受音しているので、雑音源からの音  $s_2(t)$  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、相互相関関数  $R(\quad)$  は 0 に近い値になる。一方、発話音源（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。このとき、相互相関関数  $R(\quad)$  は大きい値になる。そして、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。

40

【0053】

このように、相互相関関数計算部 3 1 は、発話音源（ユーザ）から発話があった場合には、大きい相互相関関数  $R(\quad)$  を得る。

このようなことから、音声 / 非音声判定部 3 2 は、相互相関関数  $R(\quad)$  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  $R(\quad)$  が判定用しきい値  $r_2$  より大きい場合（ $R(\quad) > r_2$ ）、音声区間と判定し、それ以外の場合（ $R(\quad) < r_2$ ）、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。そして、音声 / 非音声判定部 3 2 は、その判定結果を分離処理オン / オフ制御部 1 2 及び最適化計算オン / オフ制御部 1 3 に出力する。

【0054】

50

次に第 2 の実施形態における効果を説明する。

先ず、第 2 の実施形態では、前述した第 1 の実施形態と同様な効果を得ることができる。

さらに、第 2 の実施形態では、無指向性マイクに発話音源からの音及び雑音源からの音からなる混合音が入力され、単一指向性マイクに発話音源からの音又は雑音源からの音のいずれか一方が入力され、無指向性マイクに入力された混合音の音声信号と単一指向性マイクに入力された発話音源からの音又は雑音源からの音のいずれか一方の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、発話区間を検出している。

【0055】

これにより、無指向性マイクに発話音源からの音及び雑音源からの音からなる混合音が入力され、単一指向性マイクに発話音源からの音又は雑音源からの音のいずれか一方が入力されるように、無指向性マイク及び単一指向性マイク（第 1 及び第 2 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発話区間検出システムを構築することができる。

そして、このように精度よく発話区間を検出することができるので、分離行列を劣化させることなく、最適化することができるようになる。

【0056】

次に第 3 の実施形態を説明する。

この第 3 の実施形態は、発話区間判定部 11 を具体的な構成とした音声信号処理装置 10 であり、発話区間判定部 11 が、第 1 及び第 2 マイク 1, 2 で受音した音声信号  $x_1(t)$ ,  $x_2(t)$  のパワースペクトルに基づいて発話区間を検出するように構成されている。

図 8 は、その第 3 の実施形態における発話区間判定部 11 の構成を示す。

また、前述の第 2 の実施形態と同様、第 1 マイク 1 として単一指向性マイクを使用し、第 2 マイク 2 として無指向性マイクを使用している。そして、第 1 及び第 2 マイク 1, 2 の配置についても、前記図 6 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  $s_2(t)$  を第 1 及び第 2 マイク 1, 2 で受音し、発話音源（ユーザ）からの音  $s_1(t)$  を第 2 マイク 2 だけで受音するようにしている。

【0057】

このように配置した第 1 及び第 2 マイク 1, 2 から入力された音声信号  $x_1(t)$ ,  $x_2(t)$  はそれぞれ、前述したように、第 1 及び第 2 フレーム化部 21, 22 に入力される。そして、前述したように、第 1 フレーム化部 21 では、第 1 マイク 1 から入力された音声信号  $x_1(t)$  をフレーム化（或いはフレーム分割）し、また、第 2 フレーム化部 22 では、第 2 マイク 2 から入力される音声信号  $x_2(t)$  をフレーム化（或いはフレーム分割）する。そして、このように各フレーム化部 21, 22 で複数フレームにされた音声信号  $x_1(t)$ ,  $x_2(t)$  は発話区間判定部 11 に入力される。

【0058】

発話区間判定部 11 は、図 8 に示すように、パワースペクトラム計算部 41、パワー比計算部 42 及び音声/非音声判定部 43 を備えている。このような発話区間判定部 11 において、各フレーム化部 21, 22 で複数フレームにされた音声信号  $x_1(t)$ ,  $x_2(t)$  がパワースペクトラム計算部 41 に入力される。

パワースペクトラム計算部 41 は、フレーム単位で音声信号  $x_1(t)$ ,  $x_2(t)$  の第 1 及び第 2 パワースペクトル値  $P_{x_1}(\ )$ ,  $P_{x_2}(\ )$  を算出し、その算出した第 1 及び第 2 パワースペクトル値  $P_{x_1}(\ )$ ,  $P_{x_2}(\ )$  をパワー比計算部 42 に出力する。

パワー比計算部 42 は、下記 (2) 式により、パワースペクトラム計算部 41 からの第 1 パワースペクトル値  $P_{x_1}(\ )$  と第 2 パワースペクトル値  $P_{x_2}(\ )$  との比（以下、パワー比という。） $P(\ )$  を算出する。

【0059】

10

20

30

40

50

## 【数 2】

$$P(\omega) = G \frac{P_x(\omega)}{P_y(\omega)} \quad \dots (2)$$

## 【0060】

ここで、G 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。

パワー比計算部 42 は、このようなパワー比  $P(\omega)$  を音声 / 非音声判定部 43 に出力する。

音声 / 非音声判定部 43 は、パワー比  $P(\omega)$  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。 10

前述したように、発話音源（ユーザ）と雑音源に対して前記図 6 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  $s_2(t)$  を第 1 及び第 2 マイク 1, 2 で受音し、話者音源（ユーザ）からの音  $s_1(t)$  を第 2 マイク 2 だけで受音している。

## 【0061】

これにより、雑音源からの音  $s_2(t)$  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、このときにパワースペクトラム計算部 41 で算出される第 1 及び第 2 パワースペクトル値  $P_{x_1}(\omega)$ ,  $P_{x_2}(\omega)$  は同程度になる。一方、発話音源（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  $P_{x_1}(\omega)$  よりも第 2 パワースペクトル値  $P_{x_2}(\omega)$  の方が大きくなる。このとき、パワー比計算部 42 が算出するパワー比  $P(\omega)$  は小さくなる。 20

## 【0062】

なお、このとき、雑音源や発話音源（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  $P_{x_1}(\omega)$ ,  $P_{x_2}(\omega)$  が特に変化する。

このように、発話音源（ユーザ）から発話があった場合にはパワー比  $P(\omega)$  は小さくなることから、音声 / 非音声判定部 43 は、パワー比  $P(\omega)$  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。 30

## 【0063】

ここで、パワースペクトラム計算部 41 では、パワースペクトル値  $P_{x_1}(\omega)$ ,  $P_{x_2}(\omega)$  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  $P(\omega)$  は、各周波数帯について得ることができる。

このようなことから、パワースペクトル値  $P_{x_1}(\omega)$ ,  $P_{x_2}(\omega)$  について各周波数で得ているパワー比  $P(\omega)$  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。

## 【0064】

なお、判定対象としてパワースペクトル値  $P_{x_1}(\omega)$ ,  $P_{x_2}(\omega)$  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、発話音源（ユーザ）の特性を示す特定の周波数帯のパワー比  $P(\omega)$  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  $P(\omega)$  の平均値と判定用しきい値  $p_1$  とを比較したり、又は発話音源（ユーザ）の特性を示す特定の周波数帯のパワー比  $P(\omega)$  と雑音源の特性を示す特定の周波数帯のパワー比  $P(\omega)$  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。 40

## 【0065】

そして、音声 / 非音声判定部 43 は、パワー比  $P(\omega)$  が判定用しきい値  $p_1$  未満の場合 ( $P(\omega) < p_1$ )、音声区間と判定し、それ以外の場合 ( $P(\omega) \geq p_1$ )、 50

非音声区間と判定する。ここで、音声/非音声判定部43は、このような判定をフレーム単位で行う。そして、発話区間判定部11は、このように音声/非音声判定部43で得た音声区間(発話区間)の判定結果を分離処理オン/オフ制御部12及び最適化計算オン/オフ制御部13に出力する。

【0066】

以上のように、発話区間判定部11では、パワースペクトラム計算部41が第1及び第2フレーム化部21, 22それぞれから出力されるフレーム単位で第1及び第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42に出力する。パワー比計算部42では、パワースペクトラム計算部41から出力される第1及び第2パワースペクトル値 $P_{x_1}(\quad)$ ,  $P_{x_2}(\quad)$ について、フレーム単位でパワー比 $P(\quad)$ を算出して、算出したパワー比 $P(\quad)$ を音声/非音声判定部43に出力する。

10

【0067】

音声/非音声判定部43では、パワー比 $P(\quad)$ と判定用しきい値 $p_1$ とを比較し、パワー比 $P(\quad)$ に対応するフレームが音声区間のものか、非音声区間のものを判定する。そして、音声/非音声判定部43は、その判定結果を分離処理オン/オフ制御部12及び最適化計算オン/オフ制御部13に出力する。

そして、分離処理オン/オフ制御部12は、前述したように、発話区間判定部11からの発話区間の判定結果に基づいて、分離処理部25への観測信号 $x_1(f)$ ,  $x_2(f)$ の出力をオン及びオフ制御する。また、最適化計算オン/オフ制御部13は、前述したように、発話区間判定部11からの発話区間の判定結果に基づいて、分離行列最適化計算部26のオン及びオフを制御する。

20

【0068】

このように、第3の実施形態として、発話区間判定部11を構成することにより、前述した第1の実施形態に加えて、第2の実施形態と同様な効果を得ることができる。すなわち、無指向性マイクに発話音源からの音及び雑音源からの音からなる混合音が入力され、単一指向性マイクに発話音源からの音又は雑音源からの音のいずれか一方が入力されるように、無指向性マイク及び単一指向性マイク(第1及び第2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発話区間検出システムを構築することができる。そして、このように精度よく発話区間を検出することができるので、分離行列を劣化させることなく、最適化することができるようになる。

30

【0069】

次に第4の実施形態を説明する。

この第4の実施形態は、発話区間判定部11を具体的な構成とした音声信号処理装置10であり、第1及び第2マイク1, 2で受音した音声信号 $x_1(t)$ ,  $x_2(t)$ のクロススペクトルに基づいて発話区間を検出するように構成されている。

なお、第1及び第2マイク1, 2で受音した音声信号 $x_1(t)$ ,  $x_2(t)$ のクロススペクトルに基づいて発話区間を検出する技術については、例えば、多々良潔による「複数マイクロホンを用いた音声認識用耐雑音受音系の研究」(名古屋大学大学院工学研究科, 修士論文, 2003年3月)で開示されている。

40

【0070】

以下、このような開示技術を適用して構成した第4の実施形態における発話区間判定部11の構成を説明する。図9は、その第4の実施形態における発話区間判定部11の構成を示す。

図9に示すように、発話区間判定部11は、クロススペクトル計算部51、位相抽出処理部52、位相unwrap処理部53及び主計算部60を備えている。また、主計算部60は、周波数帯域分割部61、第1乃至第N傾き計算部 $62_1 \sim 62_N$ 、ヒストグラム等計算部63及び音声/非音声判定部64を備えている。なお、第1及び第2マイク1, 2の配置については、当該第1及び第2マイク1, 2に発話音源(ユーザ)からの音 $s_1(t)$

50

が入力されるように配置されている。

【0071】

このような発話区間判定部11において、各周波数分析部23, 24から出力された信号 $x_1(f)$ ,  $x_2(f)$ がクロススペクトル計算部51に入力される。

ここで、例えば、第1マイク1と第2マイク2といった複数のマイクで受音した音声信号を時間軸上でみた場合、受音した音声信号間に位相差が生じる。これは、音源から各マイク1, 2までの距離の違いにより、音源から各マイク1, 2までの音声信号の到達時間に差が生じた結果である。

【0072】

ここで、第1マイク1と第2マイク2とにより受音した音声信号間の遅延時間を計測し、その計測した遅延時間に基づいて位相を同相化し、その後、第1マイク1と第2マイクとでそれぞれ受音した音声信号を加算して同期加算音声を得る場合を考える。例えば、M. Omologo, P. Svaizerらの文献「Acoustic event localization using a crosspower-spectrum phase based technique”, Proc. ICASSP94, pp.274-276, (1994)」に、そのように同期加算音声を得る技術が記載されている。

10

【0073】

ここで、2つのマイク1, 2で受音した音声信号 $x_1(t)$ ,  $x_2(t)$ をフーリエ変換して得られる周波数関数を $X_1(\omega)$ ,  $X_2(\omega)$ とする。ここで、 $x_2(t)$ は、下記(3)式のように $x_1(t)$ の時間移動波形であると仮定する。

$$x_2(t) = x_1(t - t_0) \cdots (3)$$

20

このように仮定した場合、周波数関数 $X_1(\omega)$ と周波数関数 $X_2(\omega)$ との関係は下記(4)式のようになる。

$$X_2(\omega) = e^{-j\omega t_0} X_1(\omega) \cdots (4)$$

そして、この周波数関数 $X_1(\omega)$ と周波数関数 $X_2(\omega)$ とからクロススペクトル $G_{12}(\omega)$ が下記(5)式として得られる。

$$G_{12}(\omega) = X_1(\omega) X_2^*(\omega) = X_1(\omega) e^{j\omega t_0} X_1^*(\omega) = |X_1(\omega)|^2 e^{j\omega t_0} \cdots (5)$$

【0074】

ここで、クロススペクトル $G_{12}(\omega)$ の指数項はスペクトル領域のチャンネル間の時間遅れに対応する。したがって、周波数関数 $X_2$ に遅延項 $e^{j\omega t_0}$ をかけた $X_2(\omega) e^{j\omega t_0}$ は、周波数関数 $X_1$ と同相化され、これにより、 $X_1(\omega) + X_2(\omega) e^{j\omega t_0}$ の逆フーリエ変換をチャンネル同期加算音声として扱うことができるようになる。

30

【0075】

クロススペクトル計算部51では、このようなクロススペクトル $G_{12}(\omega)$ を得る。そのため、第1周波数分析部23は、第1フレーム化部21からの音声信号をフーリエ変換して前記周波数関数 $X_1(\omega)$ を算出して、その周波数関数 $X_1(\omega)$ ( $x_1(f)$ )をクロススペクトル計算部51に出力する。また、第2周波数分析部24は、第2フレーム化部22からの音声信号をフーリエ変換して周波数関数 $X_2(\omega)$ を算出して、その周波数関数 $X_2(\omega)$ ( $x_2(f)$ )をクロススペクトル計算部51に出力する。ここで、第1及び第2周波数分析部23, 24は、フレーム毎に音声信号をフーリエ変換する。

40

クロススペクトル計算部51は、第1及び第2周波数分析部23, 24からの周波数関数 $X_1(\omega)$ ,  $X_2(\omega)$ である前記信号 $x_1(f)$ ,  $x_2(f)$ に基づいて、前記(5)式によりクロススペクトル $G_{12}(\omega)$ を算出する。

【0076】

なお、図10は、1フレームについての音声信号のクロススペクトルの位相を示す。ここで、図10中(A)は自動車内で発した音声について得たクロススペクトルの位相であり、図10中(B)はオフィススペース内で発した音声について得たクロススペクトルの位相であり、図10中(C)は防音室内で発した音声について得たクロススペクトルの位相であり、図10中(D)は歩道(屋外)で発した音声について得たクロススペクトルの位相である。この図10に示すように、フレーム内で(すなわち局所的に)、音源と第1

50

マイク 1 までの距離と音源と第 2 マイク 2 までの距離との差に対応して、クロススペクトルの位相が周波数に対してほぼ一定の傾きを示すことがわかる。すなわち、音源と第 1 マイク 1 までの距離と音源と第 2 マイク 2 までの距離との差に対応して、クロススペクトルの位相成分が一定の傾きを有している。

【0077】

そして、第 1 及び第 2 マイク 1, 2 で受音した音声信号の S/N 比が高ければ、そのように傾きが一定となる傾向は顕著になる。よって、第 1 及び第 2 マイク 1, 2 により音声（発話）を受音した場合のその音声信号は S/N 比が高くなり、この場合、明らかに一定の傾きを示すものになる。

クロススペクトル計算部 5 1 は、このような特性を有するクロススペクトル  $G_{1,2}$  ( ) を位相抽出部 5 2 に出力する。 10

【0078】

位相抽出部 5 2 では、クロススペクトル計算部 5 1 からのクロススペクトル  $G_{1,2}$  ( ) から位相を抽出（検出）して、その抽出結果を位相unwrap処理部 5 3 に出力する。

位相unwrap処理部 5 3 では、位相抽出部 5 2 の位相抽出結果に基づいて、クロススペクトル  $G_{1,2}$  ( ) をunwrap処理して、主計算部 6 0 の周波数帯域分割部 6 1 に出力する。

周波数帯域分割部 6 1 は、帯域分割（セグメント分割）した位相を第 1 乃至第 N 傾き計算部 6 2<sub>1</sub> ~ 6 2<sub>N</sub> それぞれに出力する。

【0079】

ここで、音声の入力されていない非音声区間フレームと音声が入力されている音声区間フレームとで、クロススペクトルの位相成分に大きな違いがある。すなわち、音声区間フレームでは、前述したようにクロススペクトルの位相が周波数に対してほぼ一定の傾きを示すが、非音声区間フレームでは、そのようにはならない。ここで、図 1 1 を用いて説明する。 20

【0080】

図 1 1 はクロススペクトルの位相を示しており、図 1 1 中 (A) は、音声区間フレームのクロススペクトルの位相であり、図 1 1 中 (B) は、非音声区間フレームのクロススペクトルの位相である。

この図 1 1 中 (A) と図 1 1 中 (B) との比較からもわかるように、非音声区間フレームでは、クロススペクトルの位相は、周波数に対して特定のトレンドをもたない。すなわち、周波数に対してクロススペクトルの位相が一定の傾きを持つ結果とはならない。これは、ノイズの位相がランダムだからである。 30

【0081】

これに対して、音声区間フレームでは、周波数に対してクロススペクトルの位相が一定の傾きをもつようになる。そして、この傾きは、音源から各マイク 1, 2 までの距離の差に対応した大きさになる。

このように、音声の入力されていない非音声区間フレームと音声が入力されている音声区間フレームとでは、クロススペクトルの位相成分に大きな違いがある。

【0082】

このようなことから、位相の回転が生じた場合にも正確にトレンドを追従するために、周波数帯域分割部 6 1 により、位相成分を小さな周波数セグメントに分割（或いは帯域分割）し、後段の第 1 乃至第 N 傾き計算部 6 2<sub>1</sub> ~ 6 2<sub>N</sub> で、最小 2 乗法を適用することでセグメント毎に傾きを計算している。この第 1 乃至第 N 傾き計算部 6 2<sub>1</sub> ~ 6 2<sub>N</sub> はそれぞれ、算出した傾きをヒストグラム等計算部 6 3 に出力する。 40

【0083】

ここで、最小 2 乗法によりセグメント毎に傾きを求める手法は、公知の技術であり、例えば、『「信号処理」「画像処理」のための入門工学社』（高井信勝著，工学社，2000）にその技術が記載されている。

ヒストグラム等計算部 6 3 は、第 1 乃至第 N 傾き計算部 6 2<sub>1</sub> ~ 6 2<sub>N</sub> が算出した前記傾きについて、ヒストグラムを得る。

## 【 0 0 8 4 】

図 1 2 は、ヒストグラム等計算部 6 3 が得たヒストグラムであり、セグメント毎に得た傾きについてのヒストグラムを示す。すなわち、この図 1 2 は、位相の傾きの分布を示し、全セグメントに対する、各傾きのセグメント数の割合、すなわち頻度を縦軸にとっている。ここで、図 1 2 中 ( A ) は、音声区間フレームについてのヒストグラムを示し、図 1 2 中 ( B ) は、非音声区間フレームについてのヒストグラムを示す。

## 【 0 0 8 5 】

図 1 2 中 ( A ) と図 1 2 中 ( B ) との比較からわかるように、音声区間フレームでは、ヒストグラムに明らかにピーク値があり、すなわち傾きがごく狭い範囲に局在しており、これにより、ある傾きについて頻度が高くなっている。すなわち、帯域毎のそれぞれの傾きが特定の傾きに集中する傾向が強くなっている。一方、非音声区間フレームでは、ヒストグラムが平滑となり、傾きが広い範囲にわたって分布している。

10

## 【 0 0 8 6 】

このヒストグラム等計算部 6 3 は、このようなヒストグラム化して得た頻度を音声 / 非音声判定部 6 4 に出力する。なお、このヒストグラム等計算部 6 3 の処理については後で具体例を説明する。

音声 / 非音声判定部 6 4 は、ヒストグラム等計算部 6 3 からの前記頻度に基づいて、音声区間と非音声区間とを判定する。例えば、前記頻度の平均値周辺の所定の範囲に含まれる傾きの出現頻度が所定のしきい値以上の場合、音声区間と判定し、頻度が所定のしきい値未満の場合、非音声区間と判定する。

20

なお、ここでは、前段の処理がフレーム単位の処理となっているので、当該フレームが、音声区間フレーム又は非音声区間フレームのいずれかであるかを判定する。音声 / 非音声判定部 6 4 は、その判定結果を分離処理オン / オフ制御部 1 2 及び最適化計算オン / オフ制御部 1 3 に出力する。

## 【 0 0 8 7 】

次にヒストグラム等計算部 6 3 の具体的な構成を説明する。図 1 3 は、その構成例を示す。

ヒストグラム等計算部 6 3 は、第 1 乃至第 N 傾き計算部  $6 2_1 \sim 6 2_N$  が算出した前記傾きのうちから頻度が高い ( 最頻度の ) 傾きを算出する構成として、第 1 スイッチ  $6 3 S_1$ 、第 2 スイッチ  $6 3 S_2$  及び最頻値計算部  $6 3 C$  を備えている。これにより、第 1 スイッチ  $6 3 S_1$  を一定時間オン ( 閉 ) にして、第 1 乃至第 N 傾き計算部  $6 2_1 \sim 6 2_N$  が算出した一定時間の前記傾きのデータ ( 或いはデータベース )  $6 3 D_1$  を作成する。このとき、第 2 スイッチ  $6 3 S_2$  については、オフ ( 開 ) にしておく。そして、データ  $6 3 D_1$  を作成したら、第 2 スイッチ  $6 3 S_2$  をオン ( 閉 ) にして、そのデータ  $6 3 D_1$  を最頻値計算部  $6 3 C$  に出力する。

30

## 【 0 0 8 8 】

最頻値計算部  $6 3 C$  では、データ  $6 3 D_1$  から前記図 1 2 に示すような前記傾きについてのヒストグラムを作成して、そのヒストグラム中の最頻度の傾き ( 以下、最頻傾きという。 )  $\theta_0$  を算出する。

なお、最頻度の傾きを算出するようにしてもよいが、平均値の傾き  $\theta_0$  を算出したり、或いは最頻度の傾きと傾きの平均値とを組み合わせた傾き  $\theta_0$  を算出するようにしてもよい。これにより、各帯域の傾きが特定の傾きに集中する傾向が強くなったとき、当該特定の傾きの値そのもの或いはそれに近い傾きの値を得ることができる。なお、本実施の形態では、最頻値計算部  $6 3 C$  が最頻傾き  $\theta_0$  を算出しているものとする。

40

## 【 0 0 8 9 】

そして、最頻値計算部  $6 3 C$  は、算出した最頻傾き  $\theta_0$  を前記音声 / 非音声判定部 6 4 に出力する。ここで、最頻傾き  $\theta_0$  をデータ  $6 3 D_2$  として前記音声 / 非音声判定部 6 4 に出力する。

音声 / 非音声判定部 3 4 では、ヒストグラム等計算部 6 3 からの最頻傾き  $\theta_0$  に基づいて、音声区間と非音声区間とを判定する。

50

## 【0090】

なお、先の説明では、音声/非音声判定部34がヒストグラム等計算部63からの前記頻度に基づいて音声区間と非音声区間とを判定する場合について説明した。ここでは、音声/非音声判定部64は、ヒストグラム等計算部63からの最頻傾き0と第1乃至第N傾き計算部62<sub>1</sub>~62<sub>N</sub>が算出した前記傾きiに基づいて、音声区間と非音声区間とを判定しており、これに対応して、音声/非音声判定部64に、第1乃至第N傾き計算部62<sub>1</sub>~62<sub>N</sub>が算出した前記傾きが入力されるようになっている。

## 【0091】

すなわち、音声/非音声判定部64は、第1乃至第N傾き計算部62<sub>1</sub>~62<sub>N</sub>が算出した前記傾きiと最頻傾き0とを下記(6)式により比較する。

$$|i - 0| < \dots (6)$$

ここで、 $\dots$ は判定用のしきい値(傾きしきい値)である。

音声/非音声判定部34は、この(6)式の条件が満たされていることが所定の割合を超えた場合(YES)、音声区間と判定し、そうでない場合(NO)、非音声区間と判定する。そして、音声/非音声判定部64は、その判定結果を分離処理オン/オフ制御部12及び最適化計算オン/オフ制御部13に出力する。

## 【0092】

次に第4の実施形態における効果を説明する。

先ず、第4の実施形態では、前述した第1の実施形態と同様な効果を得ることができる。

さらに、第4の実施形態では、第1及び第2マイク1,2に入力された音声信号間のクロススペクトルの位相を検出し、その検出したクロススペクトルの位相の周波数に対する傾きに基づいて、当該複数のマイクロホンが受音した音声信号中の発話区間を検出している。すなわち、音声が入力(発話入力)されていない音声信号と音声が入力(発話入力)されている音声信号とをクロススペクトルでみた場合に、そのクロススペクトルの位相成分に大きな違いがあることを利用して、当該複数のマイクロホンが受音した音声信号中の発話区間を検出している。具体的には、クロススペクトルの位相を帯域分割(セグメント分割)し、帯域毎(セグメント毎)の位相の傾きからヒストグラムを生成し、そのヒストグラムから頻度(具体的には最頻値)を得て、その頻度に基づいて、発話区間を検出している。これにより、精度よく発話区間を検出することができる。そして、このように精度よく発話区間を検出することができるので、分離行列を劣化させることなく、最適化することができるようになる。

## 【0093】

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x1(t), x2(t)を、直接第1及び第2フレーム化部21,22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x1(t), x2(t)を、AD(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21,22に入力するようにする。これを、図2に示した実施形態の音声信号処理装置10の構成に適用すると、図14に示すような構成になる。

## 【0094】

この図14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x1(t), x2(t)をそれぞれ、第1及び第2AD変換部71,72でAD変換した後、第1及び第2フレーム化部21,22に入力する。

ここで、第1及び第2AD変換部71,72でAD変換されたデータ形式は、例えば11025Hz、16bit、リニアPCMである。また、第1及び第2フレーム化部21,22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。

## 【0095】

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。

10

20

30

40

50

また、前述の実施形態の説明において、発話区間判定部 11 は、検出対象音源からの音を検出する対象音件手段又は発話源からの音の発話区間を検出する発話区間検出手段を実現しており、分離処理オン/オフ制御部 12、最適化計算オン/オフ制御部 13 及び分離行列最適化計算部 26 は、前記対象音検出手段又は発話区間検出手段が検出対象音源からの音又は発話区間を検出した場合、第 1 及び第 2 マイクに入力された音信号を用いて分離行列を最適化する分離行列最適化手段を実現しており、分離処理部 25 は、前記分離行列最適化手段が最適化した分離行列を用いて、混合音から検出対象音源の音又は発話源からの音と雑音源からの音とを分離する分離手段を実現している。

#### 【0096】

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。

#### 【0097】

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、検出対象音源からの音の検出の有無により分離行列の最適化の実行を切替可能にするとともに、検出対象音源からの音を検出した場合、分離行列の最適化を行い、検出対象音源からの音を検出できない場合、分離行列の最適化を行わないように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】

#### 【0098】

【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】前記第 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。

【図 3】前記第 1 の実施形態の音声信号処理装置の一連の動作順序を示すフローチャートである。 30

【図 4】前記第 1 の実施形態における効果の説明に使用した図である。

【図 5】本発明の第 2 の実施形態における発話区間判定部の構成を示すブロック図である。

【図 6】前記第 2 の実施形態におけるマイクの配置を示す図である。

【図 7】前記第 2 の実施形態におけるマイクの他の配置を示す図である。

【図 8】本発明の第 3 の実施形態における発話区間判定部の構成を示すブロック図である。

【図 9】本発明の第 4 の実施形態における発話区間判定部の構成を示すブロック図である。 40

【図 10】各環境のクロススペクトルの位相を示す特性図である。

【図 11】クロススペクトルの位相を示す特性図であり、(A)は、音声区間フレームのクロススペクトルの位相を示す特性図であり、(B)は、非音声区間フレームのクロススペクトルの位相を示す特性図である。

【図 12】クロススペクトルの位相に基づいて得たヒストグラムを示す特性図であり、(A)は、音声区間フレームのヒストグラムを示す特性図であり、(B)は、非音声区間フレームのヒストグラムを示す特性図である。

【図 13】前記第 4 の実施形態におけるヒストグラム等計算部などの構成を示すブロック図である。

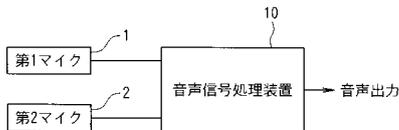
【図 14】前記第 1 の実施形態の他の構成例を示すブロック図である。 50

【符号の説明】

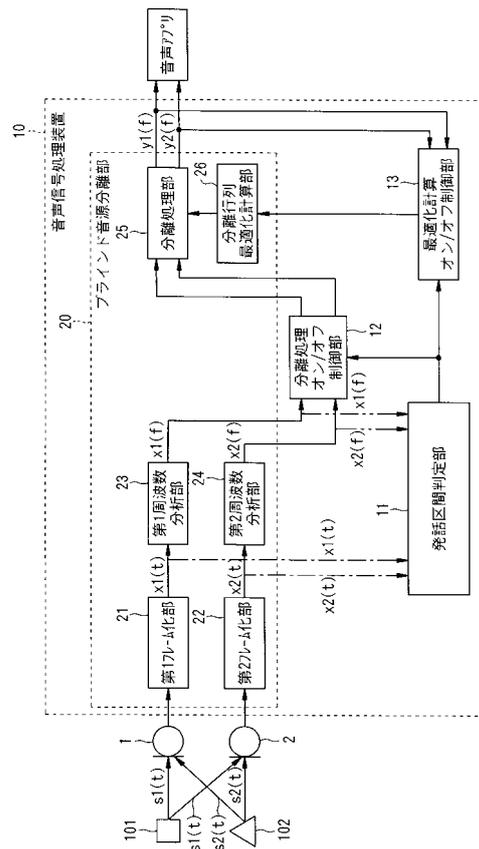
【0099】

- 1, 2    マイク
- 10    音声信号処理装置
- 11    発話区間判定部
- 12    分離処理オン/オフ制御部
- 13    最適化計算オン/オフ制御部
- 20    ブラインド音源分離部
- 21, 22   フレーム化部
- 23, 24   周波数分析部
- 25    分離処理部
- 26    分離行列最適化計算部
- 101    話者音源
- 102    雑音源

【図1】

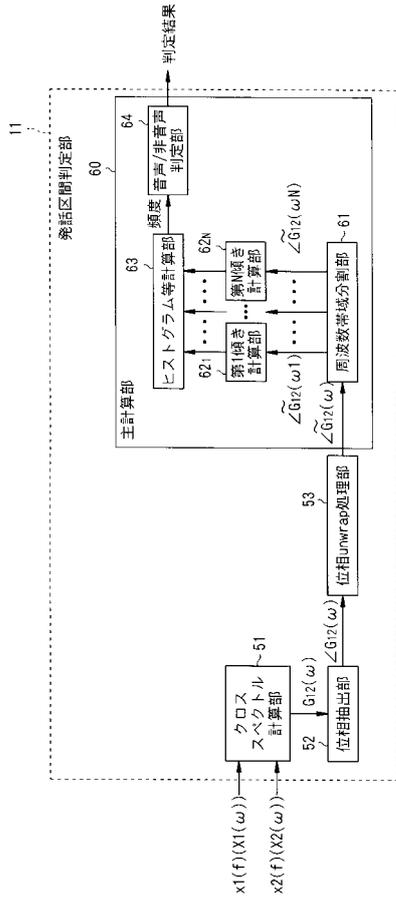


【図2】

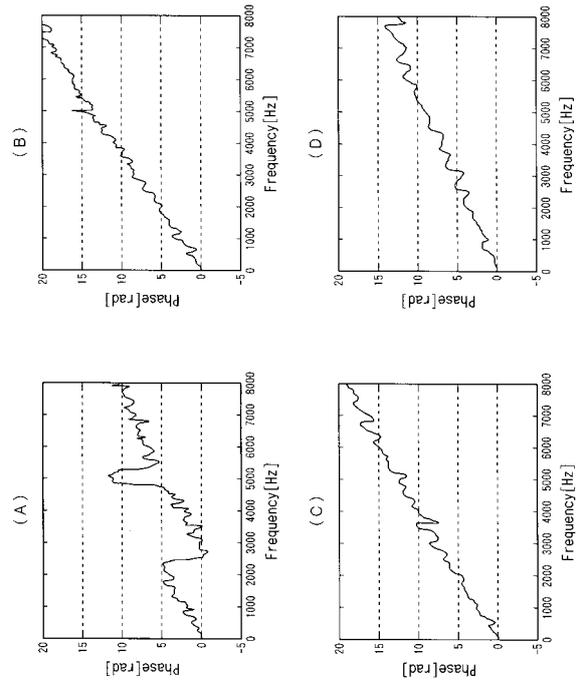




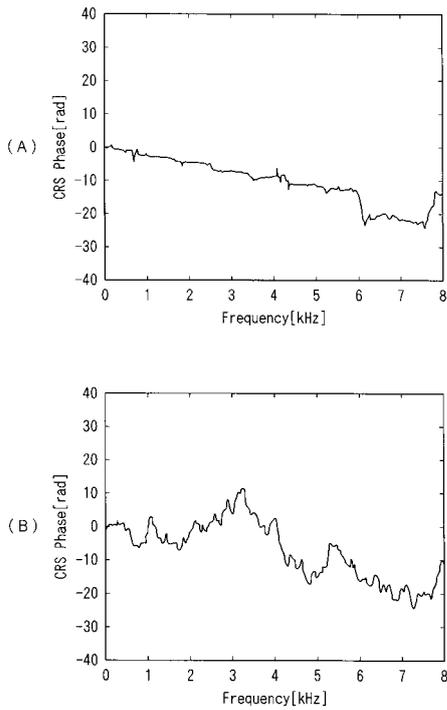
【図9】



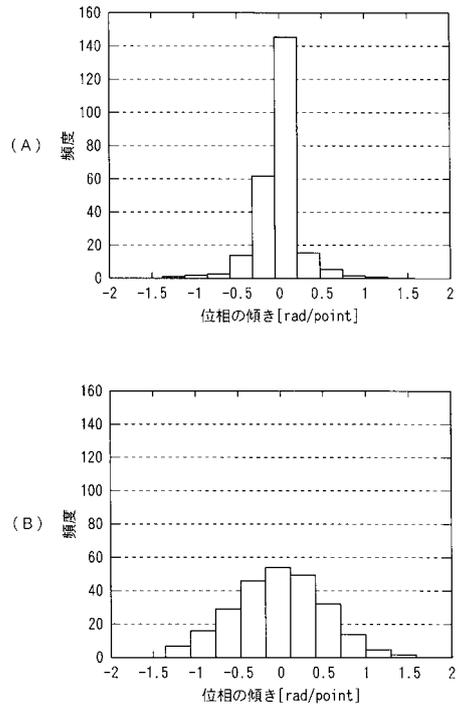
【図10】



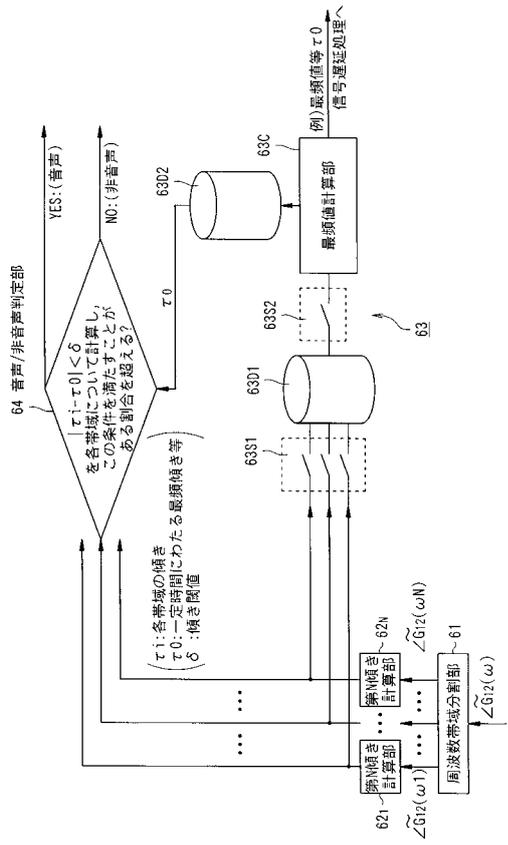
【図11】



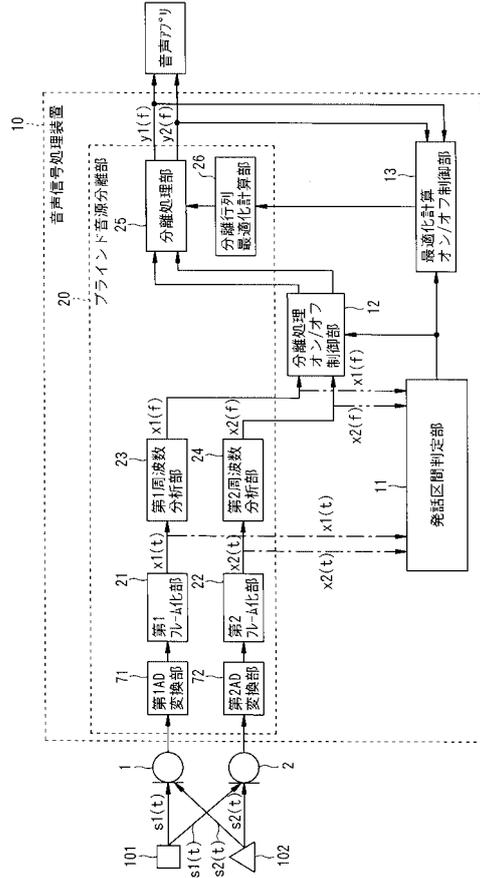
【図12】



【 図 1 3 】



【 図 1 4 】



## フロントページの続き

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

G 1 0 L 15/28  
G 1 0 L 21/02  
H 0 4 R 3/00

F I

G 1 0 L 3/02 3 0 1 F  
G 1 0 L 3/00 5 5 1 A  
G 1 0 L 9/08 3 0 1 A

テーマコード(参考)