

(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2019-101385

(P2019-101385A)

(43) 公開日 令和1年6月24日(2019.6.24)

(51) Int.Cl.	F I	テーマコード (参考)
G 1 0 L 15/04 (2013.01)	G 1 0 L 15/04 3 0 0 A	
G 1 0 L 15/00 (2013.01)	G 1 0 L 15/04 3 0 0 B	
	G 1 0 L 15/00 2 0 0 C	

審査請求 未請求 請求項の数 6 O L (全 19 頁)

(21) 出願番号	特願2017-235977 (P2017-235977)	(71) 出願人	000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号
(22) 出願日	平成29年12月8日 (2017.12.8)	(74) 代理人	100099759 弁理士 青木 篤
		(74) 代理人	100119987 弁理士 伊坪 公一
		(74) 代理人	100133835 弁理士 河野 努
		(74) 代理人	100135976 弁理士 宮本 哲夫
		(72) 発明者	鷲尾 信之 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

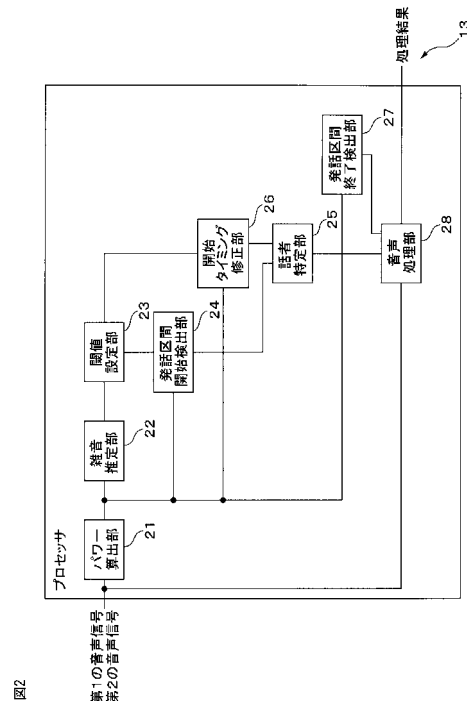
(54) 【発明の名称】 音声処理装置、音声処理方法及び音声処理用コンピュータプログラム

(57) 【要約】

【課題】 音声信号中で複数の話者の何れかが発話を開始したタイミングを誤検出しても、発話した話者に応じた処理を音声信号に適用できる音声処理装置を提供する。

【解決手段】 音声処理装置は、第1の音声入力部(11-1)により生成された第1の音声信号及び第2の音声入力部(11-2)により生成された第2の音声信号の少なくとも一方に基づいて、何れかの話者が発話を開始したタイミングを検出する発話区間開始検出部(24)と、検出された発話を開始したタイミングを修正するかどうかを判定する開始タイミング修正部(26)と、発話を開始したタイミングが修正されると、修正されたタイミング以降の第1及び第2の音声信号に基づいて、発話した話者を特定する話者特定部(25)と、特定された話者に応じた処理を、修正された発話を開始したタイミング以降の第1及び第2の音声信号の少なくとも一方に対して実行する音声処理部(28)とを有する。

【選択図】 図2



【特許請求の範囲】**【請求項 1】**

第 1 の音声入力部により生成された第 1 の音声信号及び第 2 の音声入力部により生成された第 2 の音声信号の少なくとも一方に基づいて、複数の話者の何れかが発話を開始したタイミングを検出する発話区間開始検出部と、

検出された前記発話を開始したタイミング以降における、前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方に基づいて、前記発話を開始したタイミングを修正するか否かを判定する開始タイミング修正部と、

前記発話を開始したタイミングが修正されると、修正された前記発話を開始したタイミング以降の前記第 1 の音声信号及び前記第 2 の音声信号に基づいて、前記複数の話者のうちの発話した話者を特定する話者特定部と、

特定された前記話者に応じた処理を、修正された前記発話を開始したタイミング以降の前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方に対して実行する音声処理部と、

を有する音声処理装置。

【請求項 2】

前記話者特定部は、前記発話を開始したタイミングが検出されると、当該タイミング以降の前記第 1 の音声信号及び前記第 2 の音声信号に基づいて、前記複数の話者のうちの発話した話者を特定し、

前記音声処理部は、前記発話を開始したタイミングが検出されたときに特定された前記話者に応じた第 1 の処理を、前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方に対して実行し、

前記音声処理部は、前記発話を開始したタイミングが修正されたときに前記第 1 の処理を停止する、請求項 1 に記載の音声処理装置。

【請求項 3】

前記音声処理部は、前記発話を開始したタイミングが検出されたときに特定された前記話者と、前記発話を開始したタイミングが修正されたときに特定された前記話者とが異なる場合、前記第 1 の処理を停止する、請求項 2 に記載の音声処理装置。

【請求項 4】

前記発話区間開始検出部は、前記第 1 の音声信号及び前記第 2 の音声信号のそれぞれについて、当該音声信号を分割した所定長を持つフレームごとに、前記音声信号の周期性の強さを表すピッチゲインを算出し、前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方について前記ピッチゲインが所定の閾値以上となるフレームを、前記発話を開始したタイミングとして検出し、

前記開始タイミング修正部は、前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方について、前記発話を開始したタイミングが検出されたときの前記ピッチゲインよりも前記ピッチゲインが所定のオフセット以上大きくなるフレームを検出すると、当該フレームを前記発話を開始したタイミングとして修正する、請求項 1 ~ 3 の何れか一項に記載の音声処理装置。

【請求項 5】

第 1 の音声入力部により生成された第 1 の音声信号及び第 2 の音声入力部により生成された第 2 の音声信号の少なくとも一方に基づいて、複数の話者の何れかが発話を開始したタイミングを検出し、

検出された前記発話を開始したタイミング以降における、前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方に基づいて、前記発話を開始したタイミングを修正するか否かを判定し、

前記発話を開始したタイミングが修正されると、修正された前記発話を開始したタイミング以降の前記第 1 の音声信号及び前記第 2 の音声信号に基づいて、前記複数の話者のうちの発話した話者を特定し、

特定された前記話者に応じた処理を、修正された前記発話を開始したタイミング以降の

10

20

30

40

50

前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方に対して実行する、
ことを含む音声処理方法。

【請求項 6】

第 1 の音声入力部により生成された第 1 の音声信号及び第 2 の音声入力部により生成された第 2 の音声信号の少なくとも一方に基づいて、複数の話者の何れかが発話を開始したタイミングを検出し、

検出された前記発話を開始したタイミング以降における、前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方に基づいて、前記発話を開始したタイミングを修正するか否かを判定し、

前記発話を開始したタイミングが修正されると、修正された前記発話を開始したタイミング以降の前記第 1 の音声信号及び前記第 2 の音声信号に基づいて、前記複数の話者のうちの発話した話者を特定し、

特定された前記話者に応じた処理を、修正された前記発話を開始したタイミング以降の前記第 1 の音声信号及び前記第 2 の音声信号の少なくとも一方に対して実行する、ことをコンピュータに実行させるための音声処理用コンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、例えば、話者の声を表す音声信号を処理する音声処理装置、音声処理方法及び音声処理用コンピュータプログラムに関する。

【背景技術】

【0002】

音声信号から話者が発した語句を認識し、認識した語句を他の言語に翻訳したり、認識した語句をクエリとしてネットワークまたはデータベース上で探索するといったアプリケーションが開発されている。このようなアプリケーションでは、音声信号中で話者が発話している区間が検出され、検出された区間に対してアプリケーションに応じた音声処理が実行される。

【0003】

場合によっては、複数の話者のそれぞれの声が音声処理の対象となり、かつ、話者に応じて、実行すべき処理が異なることがある。そこで、音声入力部に入力された 2 以上の使用者の音声信号を、使用者ごとに分離し、分離された使用者ごとの音声信号を認識し、その認識結果を表示部の使用者ごとに対応する表示領域に表示させる技術が提案されている（例えば、特許文献 1 を参照）。

【先行技術文献】

【特許文献】

【0004】

【特許文献 1】特開 2015 - 106014 号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、音声処理を行う装置の周囲の環境に応じて、音声信号に含まれる雑音成分の大きさが変動する。そのため、話者が発話していないにもかかわらず、音声信号中に含まれる雑音により、話者が発話を開始したタイミングが誤検出されることがある。このような場合、上記の技術では、音声信号中において実際には発話していない一方の話者の声として分離された区間中において、他方の話者が発話を開始すると、他方の話者が発話している区間も、発話していない方の話者と関連付けられてしまう。その結果として、発話している話者の声を含む区間に対して、発話していない方の話者に対する音声処理が行われてしまうことがある。

【0006】

一つの側面では、本発明は、音声信号中で複数の話者の何れかが発話を開始したタイミ

10

20

30

40

50

ングを誤検出しても、発話した話者に応じた処理を音声信号に適用できる音声処理装置を提供することを目的とする。

【課題を解決するための手段】

【0007】

一つの実施形態によれば、音声処理装置が提供される。この音声処理装置は、第1の音声入力部により生成された第1の音声信号及び第2の音声入力部により生成された第2の音声信号の少なくとも一方に基づいて、複数の話者の何れかが発話を開始したタイミングを検出する発話区間開始検出部と、検出された発話を開始したタイミング以降における、第1の音声信号及び第2の音声信号の少なくとも一方に基づいて、何れかの話者が発話を開始したタイミングを修正するか否かを判定する開始タイミング修正部と、何れかの話者が発話を開始したタイミングが修正されると、修正された発話を開始したタイミング以降の第1の音声信号及び第2の音声信号に基づいて、複数の話者のうちの発話した話者を特定する話者特定部と、特定された話者に応じた処理を、修正された発話を開始したタイミング以降の第1の音声信号及び第2の音声信号の少なくとも一方に対して実行する音声処理部とを有する。

10

【発明の効果】

【0008】

音声信号中で複数の話者の何れかが発話を開始したタイミングを誤検出しても、発話した話者に応じた処理を音声信号に適用できる。

【図面の簡単な説明】

20

【0009】

【図1】一つの実施形態による音声処理装置の概略構成図である。

【図2】音声処理に関する音声処理装置のプロセッサの機能ブロック図である。

【図3】本実施形態による話者特定の説明図である。

【図4】発話区間開始タイミングの修正についての説明図である。

【図5】話者と音声処理の対応関係の一例を示す図である。

【図6】発話区間の開始タイミングの修正と音声処理の関係の一例を示す図である。

【図7】音声処理の動作フローチャートである。

【図8】実施形態またはその変形例による音声処理装置が実装されたサーバクライアントシステムの概略構成図である。

30

【発明を実施するための形態】

【0010】

以下、図を参照しつつ、実施形態による音声処理装置について説明する。

この音声処理装置は、音声信号中で複数の話者の何れかが発話している区間（以下、単に発話区間と呼ぶ）を検出し、検出した発話区間において発話した話者を特定する。そしてこの音声処理装置は、発話区間に対して、特定した話者に応じた処理を実行する。ここで、この音声処理装置は、雑音の大きさの変動などにより発話区間の開始タイミングを誤検出した場合に備えて、発話区間の開始検出後の音声信号に基づいて、発話区間の開始を修正すべきか否かを判定する。この音声処理装置は、発話区間の開始タイミングを修正すると、修正された開始タイミングから実際の発話区間が開始されたものとして、発話した話者を再度特定する。そしてこの音声処理装置は、再特定された話者に応じた処理を、再検出された開始タイミング以降の発話区間に対して実行する。

40

【0011】

この音声処理装置は、音声信号を利用するユーザインターフェースを採用する様々な装置、例えば、ナビゲーションシステム、電話会議システム、携帯電話機またはコンピュータなどに実装できる。本実施形態では、この音声処理装置は、話者ごとに異なる言語の翻訳処理を行う多言語翻訳装置に実装されるものとする。

【0012】

図1は、一つの実施形態による音声処理装置の概略構成図である。音声処理装置1は、二つのマイクロホン11-1、11-2と、二つのアナログ/デジタルコンバータ12-

50

1、12-2と、プロセッサ13と、メモリ14と、表示装置15とを有する。なお、音声処理装置1は、さらに、スピーカ(図示せず)及び他の機器と通信するための通信インターフェース(図示せず)を有していてもよい。

【0013】

マイクロホン11-1、11-2は、それぞれ、音声入力部の一例であり、互いに対して所定の間隔を空けて設置される。例えば、マイクロホン11-1は、マイクロホン11-2よりも、複数の話者のうちの一人(便宜上、第1の話者と呼ぶ)の近くに設置される。また、マイクロホン11-2は、マイクロホン11-1よりも、複数の話者のうちの他の一人(便宜上、第2の話者と呼ぶ)の近くに設置される。そしてマイクロホン11-1、11-2は、それぞれ、複数の話者の何れかの声を含む、音声処理装置1の周囲の音を

10

【0014】

A/Dコンバータ12-1は、マイクロホン11-1から受け取ったアナログ音声信号を所定のサンプリングレートでサンプリングすることにより、その音声信号をデジタル化する。なお、サンプリングレートは、例えば、音声信号から話者の声を解析するために必要な周波数帯域がナイキスト周波数以下となるよう、例えば、16kHz~32kHzに設定される。そしてA/Dコンバータ12-1は、デジタル化された音声信号をプロセッサ13へ出力する。同様に、A/Dコンバータ12-2は、マイクロホン11-2から受け取ったアナログ音声信号を所定のサンプリングレートでサンプリングすることにより、その音声信号をデジタル化し、デジタル化された音声信号をプロセッサ13へ出力する。

20

【0015】

なお、以下では、A/Dコンバータ12-1によりデジタル化された、マイクロホン11-1からの音声信号を第1の音声信号と呼び、A/Dコンバータ12-2によりデジタル化された、マイクロホン11-2からの音声信号を第2の音声信号と呼ぶ。

【0016】

プロセッサ13は、例えば、Central Processing Unit(CPU)と、読み書き可能なメモリ回路と、その周辺回路とを有する。プロセッサ13は、数値演算回路をさらに有していてもよい。そしてプロセッサ13は、第1の音声信号及び第2の音声信号から、何れかの話者が発話している発話区間を検出するとともに、その発話区間にて発話している話者を特定する。そしてプロセッサ13は、発話区間に対して、特定した話者に対応する言語についての音声認識処理を実行し、認識された語句を、特定した話者に対応する言語以外の言語に翻訳し、その翻訳結果を表示装置15に表示させる。

30

【0017】

さらに、プロセッサ13は、一旦発話区間の開始タイミングを検出した後に、発話区間の開始タイミングを修正するか否か判定する。そして発話区間の開始タイミングが修正された場合には、プロセッサ13は、修正された発話区間の開始タイミング以降における、第1及び第2の音声信号に基づいて、発話している話者を再度特定する。プロセッサ13は、修正された開始タイミング以降の発話区間に対して、再特定された話者に対応する言語についての音声認識処理及び翻訳処理を実行する。

40

なお、音声処理の詳細については後述する。

【0018】

メモリ14は、例えば、読み書き可能な不揮発性の半導体メモリと、読み書き可能な揮発性の半導体メモリとを有する。さらに、メモリ14は、磁気記録媒体あるいは光記録媒体及びそのアクセス装置を有していてもよい。そしてメモリ14は、プロセッサ13上で実行される音声処理で利用される各種のデータ及び音声処理の途中で生成される各種のデータを記憶する。

【0019】

50

表示装置 15 は、例えば、液晶ディスプレイまたは有機 EL ディスプレイとすることができる。そして表示装置 15 は、プロセッサ 13 から受け取った表示用のデータ、例えば、何れかの話者が発話した内容、あるいは、その内容を話者が使用した言語（例えば、日本語）から他言語（例えば、英語）に翻訳して得られた文字列を表示する。

【0020】

以下、プロセッサ 13 の詳細について説明する。

【0021】

図 2 は、音声処理に関するプロセッサ 13 の機能ブロック図である。プロセッサ 13 は、パワー算出部 21 と、雑音推定部 22 と、閾値設定部 23 と、発話区間開始検出部 24 と、話者特定部 25 と、開始タイミング修正部 26 と、発話区間終了検出部 27 と、音声処理部 28 とを有する。

10

プロセッサ 13 が有するこれらの各部は、例えば、プロセッサ 13 上で動作するコンピュータプログラムにより実現される機能モジュールである。あるいは、プロセッサ 13 が有するこれらの各部は、その各部の機能を専用の回路として、プロセッサ 13 に組み込まれてもよい。

【0022】

プロセッサ 13 は、第 1 及び第 2 の音声信号のそれぞれを所定長を持つフレームを処理単位として音声処理を実行する。フレーム長は、例えば、10msec ~ 20msec に設定される。そのため、プロセッサ 13 は、第 1 及び第 2 の音声信号のそれぞれをフレームごとに分割し、各フレームをパワー算出部 21 及び音声処理部 28 へ入力する。

20

【0023】

パワー算出部 21 は、第 1 及び第 2 の音声信号のそれぞれについて、フレームが入力される度に、そのフレームについてのパワーを算出する。パワー算出部 21 は、例えば、フレームごとに、次式に従ってパワーを算出する。

【数 1】

$$Spow(k) = \sum_{n=0}^{N-1} s_k(n)^2 \quad (1)$$

30

ここで、 $s_k(n)$ は、最新のフレーム（現フレームとも呼ぶ）の n 番目のサンプリング点の信号値を表す。 k はフレーム番号である。また N は、一つのフレームに含まれるサンプリング点の総数を表す。そして $Spow(k)$ は、現フレームのパワーを表す。

【0024】

なお、パワー算出部 21 は、各フレームについて、複数の周波数のそれぞれごとにパワーを算出してもよい。この場合、パワー算出部 21 は、フレームごとに、第 1 及び第 2 の音声信号を、時間周波数変換を用いて時間領域から周波数領域のスペクトル信号に変換する。なお、パワー算出部 21 は、時間周波数変換として、例えば、高速フーリエ変換 (Fast Fourier Transform, FFT) を用いることができる。そしてパワー算出部 21 は、第 1 及び第 2 の音声信号のそれぞれについて、周波数ごとに、その周波数に含まれるスペクトル信号の 2 乗和を、その周波数のパワーとして算出できる。そしてパワー算出部 21 は、フレームごとに、人の声が含まれる周波数帯域（例えば、100Hz ~ 20kHz）に含まれる各周波数のパワーの和を、そのフレームのパワーとして算出してもよい。

40

【0025】

パワー算出部 21 は、第 1 及び第 2 の音声信号のそれぞれについて、フレームごとのパワーを、雑音推定部 22、発話区間開始検出部 24、話者特定部 25、開始タイミング修正部 26 及び発話区間終了検出部 27 へ出力する。

【0026】

雑音推定部 22 は、第 1 及び第 2 の音声信号のそれぞれについて、フレームごとに、そ

50

のフレームにおける音声信号中の推定雑音成分を算出する。本実施形態では、雑音推定部 22 は、直前のフレームにおける推定雑音成分を、現フレームのパワーを用いて次式に従って更新することで、現フレームの推定雑音成分を算出する。

【数 2】

$$Noise(k) = \beta \cdot Noise(k-1) + (1-\beta) \cdot Spow(k) \quad (2)$$

ここで、Noise(k-1) は、直前のフレームにおける推定雑音成分を表し、Noise(k) は、現フレームにおける推定雑音成分を表す。また β は、忘却係数であり、例えば、0.9 に設定される。

10

【0027】

なお、パワーが周波数ごとに算出されている場合には、雑音推定部 22 は、(2) 式に従って、推定される雑音成分を周波数ごとに算出してもよい。この場合には、(2) 式において、Noise(k-1)、Noise(k) 及び Spow(k) は、それぞれ、着目する周波数についての直前のフレームの推定雑音成分、現フレームの推定雑音成分、パワーとなる。

【0028】

雑音推定部 22 は、第 1 及び第 2 の音声信号のそれぞれについて、フレームごとの推定雑音成分を閾値設定部 23 へ出力する。

なお、後述する発話区間開始検出部 24 により、現フレームが何れかの話者の声を含む発話区間に含まれるフレームであると判定されることがある。この場合には、雑音推定部 22 は、現フレームの推定雑音成分 Noise(k) を、Noise(k-1) で置換して、(2) 式に従って再度現フレームの推定雑音成分を算出してもよい。これにより、雑音推定部 22 は、雑音成分のみを含み、信号成分を含まないと推定されるフレームに基づいて雑音成分を推定できるので、雑音成分の推定精度を向上できる。

20

【0029】

あるいは、雑音推定部 22 は、現フレームのパワーが所定の閾値以下である場合に限り、(2) 式に従って推定雑音成分を更新すればよい。そして現フレームのパワーが所定の閾値より大きい場合には、雑音推定部 22 は、Noise(k)=Noise(k-1) とすればよい。なお、所定の閾値は、例えば、Noise(k-1) に所定のオフセット値を加算した値とすることができる。

30

【0030】

閾値設定部 23 は、第 1 及び第 2 の音声信号のそれぞれについて、推定雑音成分に基づいて発話区間を検出するための閾値を設定する。例えば、閾値設定部 23 は、発話区間が検出されていない間、フレームごとに閾値を設定する。例えば、閾値設定部 23 は、第 1 の音声信号についての現フレームの推定雑音成分に、所定のオフセット値を加算した値を第 1 の音声信号に対する閾値とする。同様に、閾値設定部 23 は、第 2 の音声信号についての現フレームの推定雑音成分に、所定のオフセット値を加算した値を第 2 の音声信号に対する閾値とすればよい。

【0031】

あるいは、閾値設定部 23 は、現フレームについての、第 1 の音声信号についての推定雑音成分と第 2 の音声信号についての推定雑音成分の平均値に所定のオフセット値を加算した値を、第 1 の音声信号及び第 2 の音声信号に共通する閾値としてもよい。あるいはまた、閾値設定部 23 は、現フレームについての、第 1 の音声信号についての推定雑音成分と第 2 の音声信号についての推定雑音成分のうち大きい方に所定のオフセット値を加算した値を、第 1 の音声信号及び第 2 の音声信号に共通する閾値としてもよい。

40

【0032】

閾値設定部 23 は、第 1 及び第 2 の音声信号のそれぞれについて、発話区間の開始が検出されるまで、フレームごとに、設定した閾値を発話区間開始検出部 24 へ通知する。

【0033】

50

発話区間開始検出部 24 は、フレームごとに、そのフレームの第 1 の音声信号のパワー及び第 2 の音声信号のパワーの少なくとも一方と閾値とを比較することで、発話区間が開始されたタイミングを検出する。

【0034】

例えば、発話区間開始検出部 24 は、直前のフレームまで第 1 及び第 2 の音声信号の何れについてもパワーが対応する閾値未満であり、かつ、第 1 及び第 2 の音声信号の少なくとも一方について、現フレームのパワーが対応する閾値以上となった場合、発話区間が開始されたと判定する。そして発話区間開始検出部 24 は、現フレームを発話区間の開始タイミングとする。

【0035】

あるいは、発話区間開始検出部 24 は、フレームごとに、第 1 の音声信号及び第 2 の音声信号のうち、パワーの大きい方を、対応する閾値と比較してもよい。そして発話区間開始検出部 24 は、直前のフレームまで、パワーの大きい方が対応する閾値未満となり、かつ、現フレームにおいて、パワーの大きい方が対応する閾値以上となる場合に、現フレームを発話区間の開始タイミングとして検出してもよい。

【0036】

あるいはまた、発話区間開始検出部 24 は、第 1 の音声信号及び第 2 の音声信号の少なくとも一方について、所定数のフレームにわたって連続してパワーが対応する閾値以上となった場合、発話区間が開始されたと判定してもよい。そして発話区間開始検出部 24 は、その連続するフレームのうちの最初にパワーが閾値以上となったフレームを、発話区間の開始タイミングとして検出してもよい。

【0037】

発話区間開始検出部 24 は、発話区間が開始されたと判定すると、その旨を話者特定部 25 及び開始タイミング修正部 26 へ通知する。

【0038】

話者特定部 25 は、発話区間の開始が検出されると、その発話区間において発話している話者を特定する。例えば、話者特定部 25 は、第 1 及び第 2 の音声信号のそれぞれについて、発話区間開始検出直後の所定数（例えば、1~5）のフレームのパワーの平均値を算出する。そして話者特定部 25 は、マイクロホン 11-1、11-2 のうち、パワーの平均値が高い方の音声信号を取得したマイクロホンと対応する話者（例えば、そのマイクロホンに近い方の話者）が発話したと判定する。

【0039】

図 3 は、本実施形態による話者特定の説明図である。この例では、左から、マイクロホン 11-1、マイクロホン 11-2 の順に各マイクロホンは設置されている。そしてマイクロホン 11-1 よりも左側に第 1 の話者 301 が位置し、マイクロホン 11-2 よりも右側に第 2 の話者 302 が位置している。したがって、第 1 の話者 301 に対して、マイクロホン 11-2 よりもマイクロホン 11-1 の方が近い。そのため、第 1 の話者 301 が発話している場合、マイクロホン 11-1 により集音された第 1 の音声信号のパワーの方が、マイクロホン 11-2 により集音された第 2 の音声信号のパワーよりも大きいと推定される。したがって、発話区間開始検出直後における、第 1 の音声信号のパワーの平均値が第 2 の音声信号のパワーの平均値よりも大きい場合、第 1 の話者 301 が発話していると判定される。

【0040】

同様に、第 2 の話者 302 に対して、マイクロホン 11-1 よりもマイクロホン 11-2 の方が近い。そのため、第 2 の話者 302 が発話している場合、マイクロホン 11-2 により集音された第 2 の音声信号のパワーの方が、マイクロホン 11-1 により集音された第 1 の音声信号のパワーよりも大きいと推定される。したがって、発話区間開始検出直後における、第 2 の音声信号のパワーの平均値が第 1 の音声信号のパワーの平均値よりも大きい場合、第 2 の話者 302 が発話していると判定される。

【0041】

10

20

30

40

50

なお、話者が3名いることが想定される場合、話者特定部25は、発話区間開始検出直後の第1の音声信号のパワーの平均値と第2の音声信号のパワーの平均値との比較結果に基づいて、3名の話者のうちの何れが発話したかを判定してもよい。例えば、話者特定部25は、第1の音声信号のパワーの平均値と第2の音声信号のパワーの平均値との差の絶対値を所定のパワー差閾値と比較する。そして話者特定部25は、その差の絶対値がパワー差閾値以下である場合、マイクロホン11-1とマイクロホン11-2の並び方向に対する法線方向に位置する話者が発話したと判定してもよい。一方、話者特定部25は、その差の絶対値がパワー差閾値よりも大きく、かつ、第1の音声信号のパワーの平均値が第2の音声信号のパワーの平均値よりも大きい場合、マイクロホン11-2よりもマイクロホン11-1に近い方の話者が発話したと判定する。また、話者特定部25は、その差の絶対値がパワー差閾値よりも大きく、かつ、第2の音声信号のパワーの平均値が第1の音声信号のパワーの平均値よりも大きい場合、マイクロホン11-1よりもマイクロホン11-2に近い方の話者が発話したと判定する。

10

20

30

40

50

【0042】

あるいは、話者特定部25は、発話区間開始直後の所定数のフレームにおける第1の音声信号と第2の音声信号に基づいて、音源方向を推定し、推定した音源方向の話者が発話していると判定してもよい。この場合、話者特定部25は、例えば、発話区間開始検出直後の所定数のフレームについて、第1の音声信号と第2の音声信号間の正規化相互相関値を、互いに対する時間差をずらしながら算出する。話者特定部25は、正規化相互相関値が最も高くなる時間差を遅延時間として特定する。そして話者特定部25は、マイクロホン11-1とマイクロホン11-2間の距離と、遅延時間とに基づいて、音源方向を推定すればよい。推定された音源方向が、マイクロホン11-1とマイクロホン11-2の並び方向に対する法線方向よりも、マイクロホン11-1側を向いている場合、話者特定部25は、マイクロホン11-2よりもマイクロホン11-1に近い方の話者が発話したと判定する。なお、以下では、マイクロホン11-1とマイクロホン11-2の並び方向に対する法線方向を、マイクロホンの並び方向に対する法線方向と呼ぶ。一方、推定された音源方向が、マイクロホンの並び方向に対する法線方向よりも、マイクロホン11-2側を向いている場合、話者特定部25は、マイクロホン11-1よりもマイクロホン11-2に近い方の話者が発話したと判定する。なお、話者が3名いることが想定される場合、話者特定部25は、推定された音源方向が、マイクロホンの並び方向に対する法線方向から $\pm 45^\circ$ 未満である場合、その法線方向に位置する話者が発話したと判定してもよい。また、話者特定部25は、推定された音源方向とマイクロホンの並び方向に対する法線方向とのなす角が 45° 以上であり、かつ、その法線方向よりもマイクロホン11-1側を向いている場合、マイクロホン11-1に近い方の話者が発話したと判定する。さらに、話者特定部25は、推定された音源方向とマイクロホンの並び方向に対する法線方向とのなす角が 45° 以上であり、かつ、その法線方向よりもマイクロホン11-2側を向いている場合、マイクロホン11-2に近い方の話者が発話したと判定する。

【0043】

なお、話者特定部25は、開始タイミング修正部26により、発話区間の開始タイミングが修正されると、修正後の発話区間の開始タイミングから所定数のフレームの第1及び第2の音声信号に対して上記と同様の処理を行って、再度話者を特定する。

【0044】

話者特定部25は、特定した話者を音声処理部28へ通知する。

【0045】

開始タイミング修正部26は、発話区間開始検出部24により、発話区間の開始が検出されてからの第1及び第2の音声信号のそれぞれに基づいて、発話区間の開始タイミングを修正するか否かを判定する。

【0046】

雑音が急に大きくなることにより、発話区間開始検出部24が、雑音が急に大きくなったタイミングを、発話区間の開始タイミングと誤検出することがある。発話区間の開始タ

イミングが誤検出された後に、何れかの話者が発話を開始すると、第1及び第2の音声信号のパワーは実際の発話開始後にさらに大きくなる。そのため、誤検出された発話区間の開始タイミング直後における第1及び第2の音声信号のパワーに対して、実際の発話区間における、第1及び第2の音声信号のパワーの最大値は相対的に大きくなる。

【0047】

一方、何れかの話者が発話を継続している間、第1及び第2の音声信号には、その話者の声が含まれるので、何れかの話者が発話を継続している期間中の第1及び第2の音声信号のパワーは、そのパワーの最大値と比較してそれほど低下しない。

【0048】

そこで、開始タイミング修正部26は、発話区間の開始が検出された後の第1及び第2の音声信号のそれぞれについて、パワーの最大値を検出する。そして開始タイミング修正部26は、検出したパワーの最大値に対するパワーの低下量が所定のパワー差以上となるフレームが所定数連続すると、その連続するフレームのうちの最初のフレームを、発話区間の開始タイミングに修正する。また、開始タイミング修正部26は、第1及び第2の音声信号のそれぞれについて、発話区間検出用の閾値を、パワーの最大値から所定のパワー差を減じた値に更新する。なお、所定のパワー差は、例えば、何れかの話者が発話を継続している期間における、その話者の声による想定されるパワーの最大値とパワーの最小値の差に設定される。

【0049】

なお、開始タイミング修正部26は、発話区間の開始タイミングについての修正判定に利用する各フレームのパワーとして、パワー算出部21により算出された値そのものを用いてもよい。あるいは、開始タイミング修正部26は、その修正判定に利用する各フレームのパワーとして、パワー算出部21により算出された値から、推定雑音成分を減じた値を用いてもよい。あるいはまた、開始タイミング修正部26は、その修正判定に利用する各フレームのパワーとして、パワーの移動平均値を算出し、その移動平均値を用いてもよい。

【0050】

図4は、発話区間開始タイミングの修正についての説明図である。図4において、横軸は時間を表し、縦軸はパワーを表す。波形401は、着目する音声信号のパワーの時間変化を表す。また、波形402は、推定雑音成分のパワーの時間変化を表す。さらに、波形403は、発話区間検出用の閾値 Th の時間変化を表す。

【0051】

この例では、時刻 t_0 ~時刻 t_1 までは、着目する音声信号のパワーは閾値 Th 未満となっているので、発話区間ではないと判定される。そして時刻 t_1 の直前より、例えば、雑音が増激に大きくなることで着目する音声信号のパワーが上昇する。このとき、雑音の増加が急激であるため、閾値 Th に、雑音の増加が反映されず、その結果として、時刻 t_1 にて、着目する音声信号のパワーが閾値 Th 以上となる。そこで、時刻 t_1 にて、発話区間開始検出部24により、発話区間が開始したと判定される。

【0052】

時刻 t_1 より後の時刻 t_2 の直前にて、何れかの話者が実際に発話を開始することで、時刻 t_2 の直前より、着目する音声信号のパワーはさらに大きくなる。その結果、時刻 t_2 以降の各フレームにおいて、閾値 Th が、発話区間中のパワーの最大値 P_{max} から所定のパワー差だけ低下した値($P_{max}-$)未満となる。そのため、時刻 t_2 に、発話区間の開始タイミングが修正される。また閾値 Th は、($P_{max}-$)に更新される。その後、発話区間の開始検出後において、着目する音声信号のパワーが更新後の閾値 Th 未満となる最初のフレームの直前のフレームの時刻 t_3 において、発話区間が終了したと判定される。

【0053】

このように、閾値 Th が更新されることで、時刻 t_1 ~時刻 t_2 までの雑音しか含まれない区間が発話区間から除外されるので、発話区間が正確に求められる。

【0054】

10

20

30

40

50

なお、変形例によれば、開始タイミング修正部 2 6 は、第 1 及び第 2 の音声信号のうち、発話区間の開始検出後のパワーの最大値が大きい方の音声信号についてのみ、上記の処理を行って、発話区間の開始タイミングを修正するか否かを判定してもよい。これは、発話区間の開始検出後のパワーの最大値が大きい方の音声信号の方が、他方の音声信号よりも発話している話者の声を多く含んでいると想定されるためである。このように、一方の音声信号のみに基づいて発話区間の開始タイミングを修正するか否かを判定することで、開始タイミング修正部 2 6 は、演算量を削減できる。

【 0 0 5 5 】

開始タイミング修正部 2 6 は、発話区間の開始タイミングを修正すると、その旨を話者特定部 2 5 へ通知する。そして上記のように、話者特定部 2 5 は、発話区間の開始タイミ
10
ングを修正したことが通知されると、発話区間にて発話している話者を再度特定する。さらに、開始タイミング修正部 2 6 は、発話区間の開始タイミングを修正すると、第 1 及び第 2 の音声信号のそれぞれについて、更新された閾値 Th を発話区間終了検出部 2 7 へ通知する。

【 0 0 5 6 】

発話区間終了検出部 2 7 は、発話区間の開始が検出された以降の各フレームの第 1 及び第 2 の音声信号の少なくとも一方のパワーに基づいて発話区間が終了したか否かを判定する。

【 0 0 5 7 】

例えば、発話区間終了検出部 2 7 は、マイクロホン 1 1 - 1 と 1 1 - 2 のうち、話者特定部 2 5 により特定された話者に近い方のマイクロホンにより集音された音声信号（以下、着目音声信号と呼ぶ）のフレームのパワーを発話区間検出用の閾値と比較する。発話区間終了検出部 2 7 は、直前のフレームにおける着目音声信号のパワーが発話区間検出用の閾値以上であり、かつ、現フレームにおける着目音声信号のパワーが発話区間検出用の閾値未満である場合、直前のフレームにて発話区間が終了したと判定する。
20

【 0 0 5 8 】

あるいは、発話区間終了検出部 2 7 は、着目音声信号のパワーが発話区間検出用の閾値未満となるフレームが所定数連続すると、最初に着目音声信号のパワーが発話区間検出用の閾値未満となったフレームの直前のフレームにて発話区間が終了したと判定してもよい。
30

【 0 0 5 9 】

あるいはまた、発話区間終了検出部 2 7 は、第 1 の音声信号及び第 2 の音声信号のそれぞれについて、上記の何れかの発話区間終了検出処理を実行してもよい。そして発話区間終了検出部 2 7 は、第 1 の音声信号及び第 2 の音声信号の何れか一方、あるいは両方について、発話区間が終了したと判定される条件を満たした場合に、発話区間が終了したと判定してもよい。

【 0 0 6 0 】

なお、発話区間検出用の閾値が、開始タイミング修正部 2 6 により更新されている場合には、発話区間終了検出部 2 7 は、更新後の閾値を利用すればよい。この場合も、一旦発話区間が終了したと判定された後に、再度発話区間の開始を検出する際には、閾値設定部 2 3 により算出される、推定雑音成分に基づく閾値が利用されればよい。
40

【 0 0 6 1 】

発話区間終了検出部 2 7 は、発話区間の終了を検出すると、その旨を音声処理部 2 8 へ通知する。

【 0 0 6 2 】

音声処理部 2 8 は、発話区間の開始が検出されると、発話していると特定された話者に対応する音声処理を実行する。その際、音声処理部 2 8 は、第 1 及び第 2 の音声信号の何れに対して音声処理を実行してもよいが、例えば、マイクロホン 1 1 - 1 とマイクロホン 1 1 - 2 のうち、特定された話者に近い方のマイクロホンにて集音された音声信号に対して音声処理を実行する。発話中の話者に近い方のマイクロホンにより集音された音声信号
50

の信号対雑音比の方が、発話中の話者から遠い方のマイクロホンにより集音された音声信号の信号対雑音比よりも高いことが想定される。そのため、音声処理部 28 は、発話していると特定された話者に近い方のマイクロホンにて集音された音声信号に対して音声処理を実行することで、より適切な音声処理結果が得られる。

【0063】

図 5 は、話者と音声処理の対応関係の一例を示す図である。本実施形態では、マイクロホン 11-1 に近い方の第 1 の話者 501 が日本語を話す想定され、一方、マイクロホン 11-2 に近い方の第 2 の話者 502 が英語を話す想定される。したがって、特定された話者が第 1 の話者 501 である場合、音声処理部 28 は、第 1 の音声信号に対して日本語を対象言語とする音声認識処理を実行し、認識された発話内容に対して日本語から英語への自動翻訳処理を実行する。一方、特定された話者が第 2 の話者 502 である場合、音声処理部 28 は、第 2 の音声信号に対して英語を対象言語とする音声認識処理を実行し、認識された発話内容に対して英語から日本語への自動翻訳処理を実行する。

10

【0064】

例えば、音声処理部 28 は、発話区間中に話者が発話した内容を認識するために、処理対象となる方の音声信号の発話区間中の各フレームから、話者の声の特徴を表す複数の特徴量を抽出する。そのような特徴量として、例えば、メル周波数ケプストラムの所定の次数の係数が用いられる。そして音声処理部 28 は、例えば、各フレームの特徴量を、隠れマルコフモデルによる音響モデルに適用することで、発話区間内の音素系列を認識する。そして音声処理部 28 は、単語ごとの音素系列を表す単語辞書を参照して、発話区間の音素系列と一致する単語の組み合わせを検出することで、発話区間内の発話内容を認識する。そして音声処理部 28 は、発話内容に応じた単語の組み合わせに対して自動翻訳処理を行って、その発話内容を他言語に翻訳する。なお、音声処理部 28 は、自動翻訳処理として、様々な自動翻訳手法の何れを適用してもよい。そして音声処理部 28 は、翻訳された発話内容に応じた文字列を表示装置 15 に表示する。あるいは、音声処理部 28 は、その翻訳された文字列に音声合成処理を適用して、その文字列に対応した合成音声信号を生成し、その合成音声信号をスピーカ（図示せず）を介して再生してもよい。

20

【0065】

なお、話者が 3 名いることが想定され、特定された話者が第 1 及び第 2 の話者の何れでもない場合には、音声処理部 28 は、発話区間の第 1 及び第 2 の音声信号の何れか一方に対して、日本語及び英語以外の言語を対象とする音声認識処理を実行してもよい。あるいは、音声処理部 28 は、特定された話者が第 1 及び第 2 の話者の何れでもない場合には、前回適用された言語の音声認識処理を実行してもよい。

30

【0066】

また、音声処理部 28 は、音声処理の開始後、かつ、発話区間の終了が通知される前に、話者特定部 25 から特定された話者が再度通知され、かつ、前回通知された話者と再度通知された話者が異なっている場合には、既に開始している音声処理を停止する。そして音声処理部 28 は、再度通知された話者に対応する音声処理を実行する。これにより、発話区間の開始タイミングが誤検出されることで、特定された話者が誤っている場合でも、誤って特定された話者に対応する音声処理が継続することが防止される。

40

【0067】

図 6 は、発話区間の開始タイミングの修正と音声処理の関係の一例を示す図である。図 6 において、横軸は時間を表す。波形 601 は、第 1 及び第 2 の音声信号の一方の波形の一例である。この例では、時刻 t_1 ~ 時刻 t_2 までは、音声信号には雑音成分のみが含まれ、話者の声は含まれないものとする。一方、時刻 t_2 ~ 時刻 t_3 において、マイクロホン 11-2 に近い方の話者が発話しているとする。

【0068】

時刻 t_1 において、発話区間の開始が誤検出され、マイクロホン 11-1 に近い方の第 1 の話者が発話していると判定されたとする。この場合、誤検出された区間 602 において、音声処理部 28 により、日本語を認識対象とする音声認識処理が実行されることになる

50

。そして発話区間の開始が修正されなければ、実際に発話が始まった時刻 t_2 以降においても、日本語を認識対象とする音声認識処理が継続されるため、話者の発話内容が正確に認識されない。

【0069】

一方、本実施形態では、時刻 t_2 にて発話区間の開始タイミングが修正され、修正された発話区間の開始タイミングにて再度発話している話者が特定される。そのため、実際の発話区間603では、実際に発話している、マイクロホン11-2に近い方の第2の話者に対応する、英語を認識対象とする音声認識処理が実行される。したがって、音声処理部28は、実際に発話している話者の発話内容を正確に認識することができる。

なお、誤検出された区間に対する日本語を認識対象とする音声認識処理は、修正された発話区間の開始タイミングにて停止される。

【0070】

図7は、本実施形態による、音声処理の動作フローチャートである。プロセッサ13は、フレームごとに、下記の動作フローチャートに従って音声処理を実行する。

【0071】

パワー算出部21は、第1及び第2の音声信号のそれぞれについて、現フレームのパワー P を算出する(ステップS101)。雑音推定部22は、第1及び第2の音声信号のそれぞれについて、現フレームのパワー P と、直前のフレームにおける推定雑音成分とに基づいて、現フレームの推定雑音成分を算出する(ステップS102)。

【0072】

閾値設定部23は、直前のフレームが発話区間内か否か判定する(ステップS103)。直前のフレームが発話区間外であれば(ステップS103-No)、第1及び第2の音声信号のそれぞれについて、閾値設定部23は、推定雑音成分に基づいて閾値 Th を設定する(ステップS104)。そして発話区間開始検出部24は、第1及び第2の音声信号のそれぞれについて、現フレームのパワー P が閾値 Th 以上か否か判定する(ステップS105)。

【0073】

第1及び第2の音声信号の両方について、現フレームのパワー P が閾値 Th 未満であれば(ステップS105-No)、発話区間開始検出部24は、現フレームが発話区間に含まれないと判定する。そしてプロセッサ13は、音声処理を終了する。一方、第1及び第2の音声信号の少なくとも一方について、現フレームのパワー P が閾値 Th 以上であれば(ステップS105-Yes)、発話区間開始検出部24は、現フレームから発話区間が開始したと判定する(ステップS106)。そして発話区間開始検出部24は、現フレームを発話区間の開始タイミングとして検出する。話者特定部25は、開始した発話区間において発話した話者を、第1及び第2の音声信号に基づいて特定する(ステップS107)。さらに、音声処理部28は、第1及び第2の音声信号の何れかに対して、特定された話者に応じた処理を実行する(ステップS108)。その後、プロセッサ13は、現フレームにおける音声処理を終了する。

【0074】

また、ステップS103において、直前のフレームが発話区間に含まれる場合(ステップS103-Yes)、発話区間の開始タイミングが既に検出されている。そこで開始タイミング修正部26は、第1及び第2の音声信号のそれぞれについて、発話区間開始後のパワーの最大値 P_{max} から所定のパワー差だけ減じた値よりも閾値 Th が小さいフレームが所定数以上継続したか否か判定する(ステップS109)。

【0075】

第1及び第2の音声信号の少なくとも一方について、現フレームにおいて、 $(P_{max} -) > Th$ となるフレームが所定数以上連続している場合(ステップS109-Yes)、開始タイミング修正部26は、閾値 Th を $(P_{max} -)$ に更新する。そして開始タイミング修正部26は、発話区間の開始タイミングをその連続するフレームのうちの最初のフレームに修正する(ステップS110)。その後、プロセッサ13は、ステップS107以降の処理を

10

20

30

40

50

実行する。なお、この場合、ステップ S 1 0 8 において、音声処理部 2 8 は、発話区間の開始タイミングの修正前後において、特定された話者が異なる場合、発話区間の開始タイミングの修正前に行っている音声処理を停止する。

【 0 0 7 6 】

一方、第 1 及び第 2 の音声信号の両方について、現フレームにおいて、 $(P_{max} -) > Th$ となるフレームの継続数が所定数未満であれば (ステップ S 1 0 9 - N o)、開始タイミング修正部 2 6 は、発話区間の開始タイミングを修正しない。一方、発話区間終了検出部 2 7 は、第 1 及び第 2 の音声信号のうち、音声処理部 2 8 の音声処理対象となる方の音声信号の現フレームのパワー P が閾値 Th 未満か否か判定する (ステップ S 1 1 1)。パワー P が閾値 Th 未満であれば (ステップ S 1 1 1 - Y e s)、発話区間終了検出部 2 7 は、直前のフレームにて発話区間が終了したと判定する (ステップ S 1 1 2)。そしてプロセッサ 1 3 は、音声処理部 2 8 に対して発話区間の終了を通知する。一方、パワー P が閾値 Th 以上であれば (ステップ S 1 1 1 - N o)、発話区間終了検出部 2 7 は、現フレームも発話区間内に含まれると判定する。そしてプロセッサ 1 3 は、ステップ S 1 0 8 の処理を実行する。

10

【 0 0 7 7 】

以上に説明してきたように、この音声処理装置は、発話区間の開始が検出されると、その発話区間にて発話した話者を特定し、特定した話者に応じた音声処理を、第 1 及び第 2 の音声信号の少なくとも一方に対して実行する。そしてこの音声処理装置は、一旦発話区間の開始が検出された後に、発話区間の開始タイミングが修正されると、修正された開始タイミングにて、複数の話者のうち、その発話区間にて発話した話者を再度特定する。そして音声処理装置は、再特定された話者に応じた音声処理を第 1 及び第 2 の音声信号の少なくとも一方に対して実行する。そのため、この音声処理装置は、各音声信号中で複数の話者の何れかが発話を開始したタイミングを誤検出しても、発話した話者に応じた処理を音声信号に適用できる。

20

【 0 0 7 8 】

なお、変形例によれば、音声処理部 2 8 は、音声認識処理及び自動翻訳処理以外の処理を実行してもよい。例えば、第 1 の話者の周囲がエコーを生じ易く、第 2 の話者の周囲に雑音源が存在するとする。この場合、第 1 の話者が発話していると判定された場合、音声処理部 2 8 は、発話区間における第 1 及び第 2 の音声信号の少なくとも一方に対してエコー除去処理を実行してもよい。一方、第 2 の話者が発話していると判定された場合、音声処理部 2 8 は、発話区間における第 1 及び第 2 の音声信号の少なくとも一方に対して雑音除去処理を実行してもよい。

30

【 0 0 7 9 】

また、発話区間開始検出部 2 4 及び開始タイミング修正部 2 6 は、各フレームのパワー以外の、音声信号に含まれる、話者の声を表す特徴量に基づいて、発話区間の開始タイミングの検出、及びその開始タイミングの修正判定を行ってもよい。例えば、発話区間開始検出部 2 4 は、第 1 及び第 2 の音声信号の各フレームから、音の周期性の強さを表すピッチゲインを算出する。そして発話区間開始検出部 2 4 は、第 1 及び第 2 の音声信号の少なくとも一方について、直前のフレームのピッチゲインが閾値未満となり、かつ、現フレームのピッチゲインが閾値以上となる場合、発話区間の開始を検出してもよい。なお、ピッチゲイン g_{pitch} は、例えば、次式に従って算出される。

40

【数 3】

$$g_{pitch} = \frac{C(d_{max})}{\sum_{n=0}^{N-1} s_k(n) \cdot s_k(n)} \quad (3)$$

$$C(d) = \sum_{n=0}^{N-1} s_k(n) \cdot s_k(n-d) \quad (d = d_{low}, \dots, d_{high})$$

10

ここで、 $C(d)$ は、着目する音声信号の長期自己相関である。また $d = \{d_{low}, \dots, d_{high}\}$ は、遅延量である。そして $s_k(n)$ は、現フレーム k の n 番目の信号値である。また N は、フレームに含まれるサンプリング点の総数を表す。なお、 $(n-d)$ が負となる場合、直前のフレームの対応する信号値（すなわち、フレーム区間の重複がない場合、 $s_{k-1}(N-(n-d))$ ）が $s_k(n-d)$ として用いられる。そして遅延量 d の範囲 $\{d_{low}, \dots, d_{high}\}$ は、人の声の基本周波数(100~300Hz)に相当する遅延量が含まれるように設定される。ピッチゲインは、基本周波数において最も高くなるためである。例えば、サンプリングレートが16kHzである場合、 $d_{low}=40$ 、 $d_{high}=286$ に設定される。さらに、 d_{max} は、長期自己相関 $C(d)$ の最大値 $C(d_{max})$ に対応する遅延量であり、この遅延量はピッチ周期に相当する。

20

【0080】

一般に、ピッチゲインは、発話が始まった直後において最も大きく、発話が継続するにつれて小さくなる。そこで、開始タイミング修正部26は、第1及び第2の音声信号の少なくとも一方について、発話区間の開始を検出された直後の所定数のフレームのピッチゲインの最大値と、発話区間の開始検出後の各フレームのピッチゲインを比較する。そして開始タイミング修正部26は、そのピッチゲインの最大値よりも所定のオフセット値以上ピッチゲインが大きくなるフレームを検出すると、そのフレームに発話区間の開始タイミングを修正すればよい。

【0081】

なお、この変形例の場合、発話区間終了検出部27は、発話区間の開始検出後において、第1及び第2の音声信号の両方についてピッチゲインが閾値未満となる最初のフレームにおいて、発話区間が終了したと判定してもよい。あるいは、発話区間終了検出部27は、第1及び第2の音声信号の両方についてピッチゲインが閾値未満となるフレームが所定数連続する場合、ピッチゲインが閾値未満となった最初のフレームにおいて、発話区間が終了したと判定してもよい。また、発話区間終了検出部27は、パワーとピッチゲインが共に閾値未満となる最初のフレームにおいて、発話区間が終了したと判定してもよい。

30

【0082】

上記の実施形態または変形例による音声処理装置は、サーバクライアント型のシステムに実装されてもよい。

図8は、上記の実施形態またはその変形例による音声処理装置が実装されたサーバクライアントシステムの概略構成図である。

40

サーバクライアントシステム100は、端末110とサーバ120とを有し、端末110とサーバ120とは、通信ネットワーク130を介して互いに通信可能となっている。なお、サーバクライアントシステム100が有する端末110は複数存在してもよい。同様に、サーバクライアントシステム100が有するサーバ120は複数存在してもよい。

【0083】

端末110は、二つのマイクロホン111-1、111-2と、メモリ112と、通信インターフェース113と、プロセッサ114と、表示装置115とを有する。マイクロホン111、メモリ112及び通信インターフェース113は、例えば、プロセッサ114とバスを介して接続されている。

【0084】

50

マイクロホン 111-1、111-2は、それぞれ、音声入力部の一例である。マイクロホン 111-1は、アナログ信号である第1の音声信号を取得し、第1の音声信号をA/Dコンバータ(図示せず)へ出力する。A/Dコンバータによりデジタル化された第1の音声信号はプロセッサ114へ出力される。同様に、マイクロホン 111-2は、アナログ信号である第2の音声信号を取得し、第2の音声信号をA/Dコンバータ(図示せず)へ出力する。A/Dコンバータによりデジタル化された第2の音声信号はプロセッサ114へ出力される。

【0085】

メモリ112は、例えば、不揮発性の半導体メモリ及び揮発性の半導体メモリを有する。そしてメモリ112は、端末110を制御するためのコンピュータプログラム、端末110の識別情報、発話区間検出処理で利用される各種のデータ及びコンピュータプログラムなどを記憶する。

10

【0086】

通信インターフェース113は、端末110を通信ネットワーク130に接続するためのインターフェース回路を有する。そして通信インターフェース113は、プロセッサ114から受け取った音声信号を、端末110の識別情報とともに通信ネットワーク130を介してサーバ120へ送信する。

【0087】

プロセッサ114は、CPUとその周辺回路を有する。そしてプロセッサ114は、第1及び第2の音声信号を、端末110の識別情報とともに、通信インターフェース113及び通信ネットワーク130を介してサーバ120へ送信する。またプロセッサ114は、サーバ120から受け取った、各音声信号に対する処理結果を表示装置115に表示するか、あるいは、その処理結果に対応する合成音声信号をスピーカ(図示せず)を介して再生する。

20

【0088】

表示装置115は、例えば、液晶ディスプレイまたは有機ELディスプレイであり、各音声信号に対する処理結果を表示する。

【0089】

サーバ120は、通信インターフェース121と、メモリ122と、プロセッサ123とを有する。通信インターフェース121及びメモリ122は、プロセッサ123とバスを介して接続されている。

30

【0090】

通信インターフェース121は、サーバ120を通信ネットワーク130に接続するためのインターフェース回路を有する。そして通信インターフェース121は、第1及び第2の音声信号と端末110の識別情報とを端末110から通信ネットワーク130を介して受信してプロセッサ123に渡す。

【0091】

メモリ122は、例えば、不揮発性の半導体メモリ及び揮発性の半導体メモリを有する。そしてメモリ122は、サーバ120を制御するためのコンピュータプログラムなどを記憶する。またメモリ122は、音声処理を実行するためのコンピュータプログラム及び各端末から受信した各音声信号を記憶してもよい。

40

【0092】

プロセッサ123は、CPUとその周辺回路を有する。そしてプロセッサ123は、上記の実施形態または変形例による音声処理装置のプロセッサの各機能を実現する。そしてプロセッサ123は、受信した第1及び第2の音声信号に対する音声処理結果を通信インターフェース121及び通信ネットワーク130を介して端末110へ送信する。

【0093】

なお、端末110のプロセッサ114が、上記の実施形態または変形例による音声処理装置のプロセッサの各機能のうち、音声処理部28以外の処理を実行してもよい。この場合、端末110は、発話区間中の第1及び第2の音声信号の少なくとも何れかと、特定

50

された話者を表す情報とをサーバ120へ送信すればよい。また、端末110は、発話区間の開始タイミングを修正した場合には、修正された発話区間の開始タイミング及び再特定された話者を表す情報をサーバ120へ送信する。そしてサーバ120のプロセッサ123は、受信した第1及び第2の音声信号の少なくとも一方に対して、音声処理部28の処理を実行すればよい。

【0094】

上記の実施形態または変形例による発話区間検出装置のプロセッサが有する各機能をコンピュータに実現させるコンピュータプログラムは、磁気記録媒体または光記録媒体といったコンピュータによって読み取り可能な媒体に記録された形で提供されてもよい。

【0095】

ここに挙げられた全ての例及び特定の用語は、読者が、本発明及び当該技術の促進に対する本発明者により寄与された概念を理解することを助ける、教示的な目的において意図されたものであり、本発明の優位性及び劣等性を示すことに関する、本明細書の如何なる例の構成、そのような特定の挙げられた例及び条件に限定しないように解釈されるべきものである。本発明の実施形態は詳細に説明されているが、本発明の精神及び範囲から外れることなく、様々な変更、置換及び修正をこれに加えることが可能であることを理解されたい。

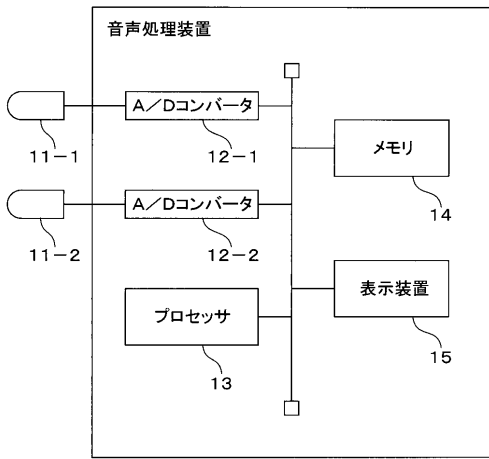
【符号の説明】

【0096】

1	音声処理装置	20
11-1、11-2	マイクロホン	
12-1、12-2	アナログ/デジタルコンバータ	
13	プロセッサ	
14	メモリ	
15	表示装置	
21	パワー算出部	
22	雑音推定部	
23	閾値設定部	
24	発話区間開始検出部	
25	話者特定部	30
26	開始タイミング修正部	
27	発話区間終了検出部	
28	音声処理部	
100	サーバクライアントシステム	
110	端末	
111-1、111-2	マイクロホン	
112	メモリ	
113	通信インターフェース	
114	プロセッサ	
115	表示装置	40
120	サーバ	
121	通信インターフェース	
122	メモリ	
123	プロセッサ	
130	通信ネットワーク	

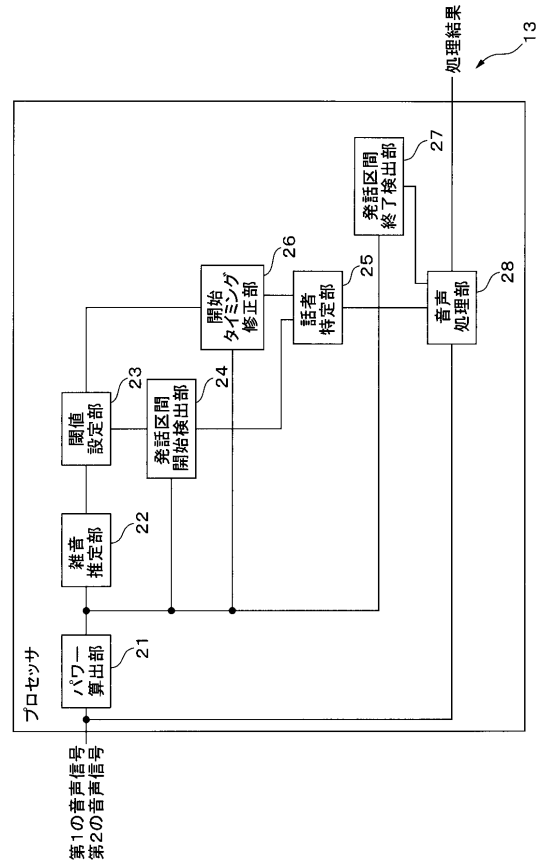
【 図 1 】

図1



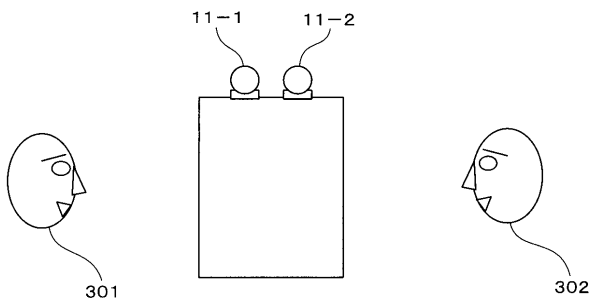
【 図 2 】

図2



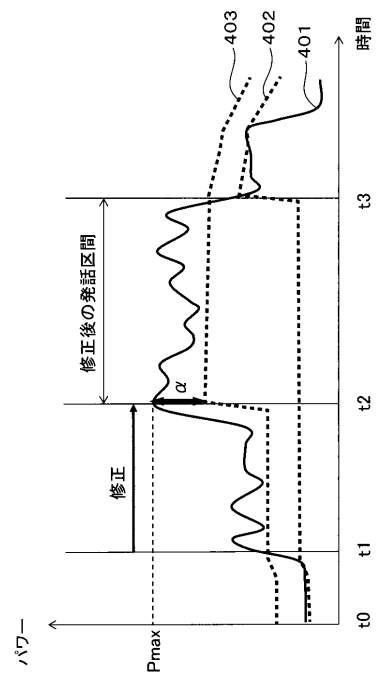
【 図 3 】

図3



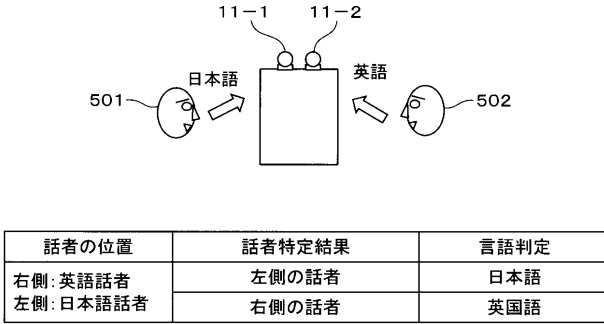
【 図 4 】

図4



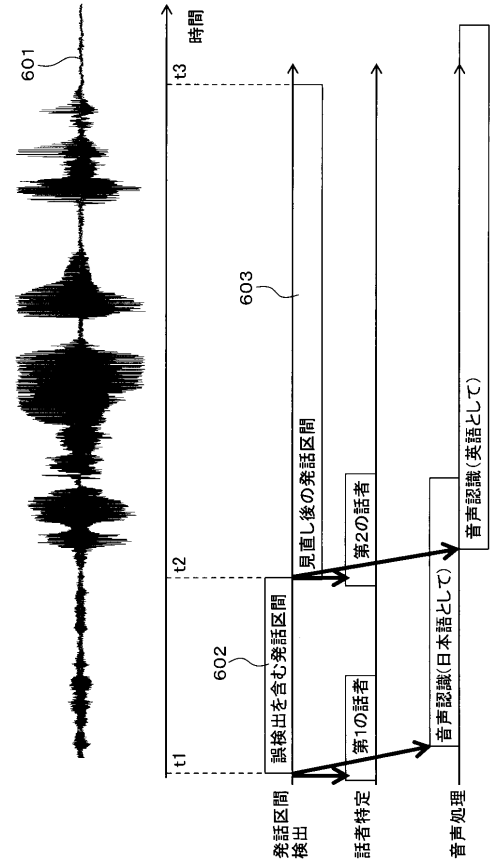
【 図 5 】

図5



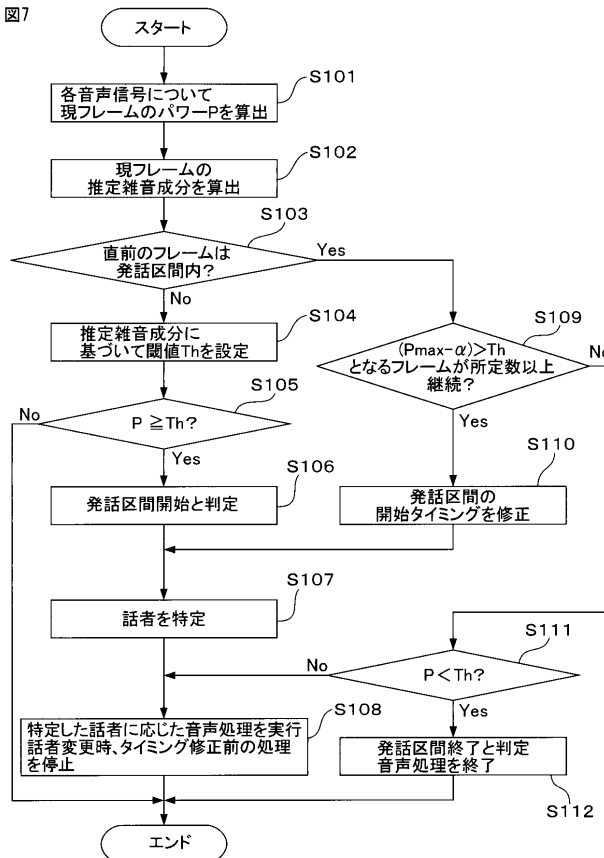
【 図 6 】

図6



【 図 7 】

図7



【 図 8 】

図8

