

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6641832号
(P6641832)

(45) 発行日 令和2年2月5日(2020.2.5)

(24) 登録日 令和2年1月8日(2020.1.8)

(51) Int.Cl.		F I	
G 1 0 L 25/51	(2013.01)	G 1 0 L	25/51
G 1 0 L 25/06	(2013.01)	G 1 0 L	25/06
G 1 0 L 15/10	(2006.01)	G 1 0 L	15/10 5 0 0 Z

請求項の数 13 (全 31 頁)

<p>(21) 出願番号 特願2015-186617 (P2015-186617)</p> <p>(22) 出願日 平成27年9月24日 (2015.9.24)</p> <p>(65) 公開番号 特開2017-62307 (P2017-62307A)</p> <p>(43) 公開日 平成29年3月30日 (2017.3.30)</p> <p>審査請求日 平成30年6月8日 (2018.6.8)</p>	<p>(73) 特許権者 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</p> <p>(74) 代理人 100147164 弁理士 向山 直樹</p> <p>(72) 発明者 外川 太郎 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</p> <p>(72) 発明者 香村 紗友梨 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</p> <p>(72) 発明者 大谷 猛 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</p> <p style="text-align: right;">最終頁に続く</p>
--	--

(54) 【発明の名称】 音声処理装置、音声処理方法および音声処理プログラム

(57) 【特許請求の範囲】

【請求項1】

第1音声が含まれる第1入力信号と、第2音声が含まれる第2入力信号を取得する取得部と、

前記第1入力信号の第1信号強度と、前記第2入力信号の第2信号強度を検出する検出部と、

前記第1信号強度の時系列と前記第2信号強度の時系列との相関係数を算出する算出部と、

前記相関係数に基づいて、前記第1音声と前記第2音声が会話している状態か否かを判定する判定部と、を有し、

前記検出部は、

前記第1入力信号に含まれる第1発話区間を前記第1信号強度に基づいて検出し、

前記第2入力信号に含まれる第2発話区間を前記第2信号強度に基づいて検出し、

前記第1発話区間と前記第2発話区間が重複する重複発話区間を検出し、

前記算出部は、

前記重複発話区間が所定の第2閾値未満の前記第1信号強度と前記第2信号強度以外の前記第1信号強度の時系列と前記第2信号強度の時系列との前記相関係数を算出することを特徴とする音声処理装置。

【請求項2】

前記判定部は、前記相関係数が負であり、前記相関係数が所定の第1閾値未満の場合に

、前記第1音声と前記第2音声が会話している状態と判定することを特徴とする請求項1記載の音声処理装置。

【請求項3】

前記第1信号強度は前記第1音声のパワーまたは信号対雑音比であり、前記第2信号強度は前記第2音声のパワーまたは信号対雑音比であることを特徴とする請求項1または請求項2記載の音声処理装置。

【請求項4】

前記検出部は、
 前記第1入力信号に含まれる第1無音区間を前記第1信号強度に基づいて検出し、
 前記第2入力信号に含まれる第2無音区間を前記第2信号強度に基づいて検出し、
 前記第1無音区間と前記第2無音区間が重複する重複無音区間を検出し、
 前記算出部は、
 前記重複無音区間が所定の第3閾値未満の前記第1信号強度と前記第2信号強度を、前記相関係数の算出に用いないことを特徴とする請求項1ないし請求項3の何れか一項に記載の音声処理装置。

10

【請求項5】

前記算出部は、前記第1信号強度の第1位相または、前記第2信号強度の第2位相を所定の範囲で変化させて複数の前記相関係数を算出し、
 前記判定部は、複数の前記相関係数のうち最小値となる前記相関係数に基づいて、前記第1音声と前記第2音声の会話している状態が否かを判定することを特徴とする請求項1ないし請求項4の何れか一項に記載の音声処理装置。

20

【請求項6】

前記検出部は、
 前記第1信号強度の時系列と前記第2信号強度の時系列との大小関係を比較し、
 前記大小関係が反転する反転回数が所定の第4閾値以上となる相関係数算出区間を検出し、
 前記算出部は、
 前記相関係数算出区間における前記第1信号強度の時系列と前記第2信号強度の時系列との前記相関係数を算出することを特徴とする請求項1ないし請求項5の何れか一項に記載の音声処理装置。

30

【請求項7】

前記取得部は、第3音声が含まれる第3入力信号を更に取得し、
 前記検出部は、前記第3入力信号の第3信号強度を更に検出し、
 前記算出部は、前記第1信号強度の時系列、前記第2信号強度の時系列または前記第3信号強度の時系列のうち、2つの信号強度の組み合わせに対する複数の前記相関係数を算出し、
 前記判定部は、複数の前記相関係数のうち、前記相関係数が最小値となる2つの前記信号強度の組み合わせに基づいて、前記第1音声、前記第2音声または前記第3音声から、会話している音声の組み合わせを判定することを特徴とする請求項1ないし請求項6の何れか一項に記載の音声処理装置。

40

【請求項8】

前記算出部は、前記最小値となる前記相関係数の算出に用いた2つの前記信号強度を加算した加算信号強度を算出し、
 前記加算信号強度の時系列と、前記最小値となる前記相関係数の算出に用いた2つの前記信号強度以外の1つの信号強度の時系列との参照相関係数を算出し、
 前記判定部は、前記相関係数が負であり、前記相関係数が所定の第1閾値未満の場合に、前記第1音声と前記第2音声の会話している状態と判定し、
 前記最小値となる前記相関係数が前記参照相関係数を下回る場合、または、前記第1閾値未満の場合に、前記参照相関係数の算出に用いた3つの前記信号強度に基づいて、会話している音声の組み合わせを判定することを特徴とする請求項7記載の音声処理装置。

50

【請求項 9】

前記検出部は、

前記第 1 入力信号に含まれる第 1 無音区間を前記第 1 信号強度に基づいて検出し、

前記第 2 入力信号に含まれる第 2 無音区間を前記第 2 信号強度に基づいて検出し、

前記第 1 無音区間と前記第 2 無音区間が重複する重複無音区間を検出し、

前記算出部は、前記最小値となる前記相関係数の算出に用いた 2 つの前記信号強度の前記重複無音区間が第 5 閾値以上の場合に、前記参照相関係数を算出することを特徴とする請求項 8 記載の音声処理装置。

【請求項 10】

前記検出部は、

前記第 1 入力信号に含まれる前記第 2 信号強度または、前記第 2 入力信号に含まれる前記第 1 信号強度を更に検出し、

前記算出部は、

前記第 1 入力信号に含まれる前記第 2 信号強度の時系列と前記第 2 入力信号に含まれる前記第 2 信号強度の時系列との第 2 相関係数、または、

前記第 2 入力信号に含まれる前記第 1 信号強度の時系列と前記第 1 入力信号に含まれる前記第 1 信号強度の時系列との第 3 相関係数を更に算出し、

前記判定部は、前記第 2 相関係数または前記第 3 相関係数に基づいて、前記第 1 音声と前記第 2 音声が会話している状態か否かを判定することを特徴とする請求項 1 ないし請求項 9 の何れか一項に記載の音声処理装置。

【請求項 11】

前記判定部は、前記相関係数が正であり、前記第 2 相関係数が所定の第 6 閾値以上の場合に、前記第 1 音声と前記第 2 音声の前記会話している状態と判定することを特徴とする請求項 10 記載の音声処理装置。

【請求項 12】

第 1 音声が含まれる第 1 入力信号と、第 2 音声が含まれる第 2 入力信号を取得し、

前記第 1 入力信号の第 1 信号強度と、前記第 2 入力信号の第 2 信号強度を検出し、

前記第 1 入力信号に含まれる第 1 発話区間を前記第 1 信号強度に基づいて検出し、

前記第 2 入力信号に含まれる第 2 発話区間を前記第 2 信号強度に基づいて検出し、

前記第 1 発話区間と前記第 2 発話区間が重複する重複発話区間を検出し、

前記重複発話区間が所定の第 2 閾値未満の前記第 1 信号強度と前記第 2 信号強度以外の前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との相関係数を算出し、

前記相関係数に基づいて、前記第 1 音声と前記第 2 音声がか会話している状態か否かを判定することを含む音声処理方法。

【請求項 13】

コンピュータに

第 1 音声が含まれる第 1 入力信号と、第 2 音声が含まれる第 2 入力信号を取得し、

前記第 1 入力信号の第 1 信号強度と、前記第 2 入力信号の第 2 信号強度を検出し、

前記第 1 入力信号に含まれる第 1 発話区間を前記第 1 信号強度に基づいて検出し、

前記第 2 入力信号に含まれる第 2 発話区間を前記第 2 信号強度に基づいて検出し、

前記第 1 発話区間と前記第 2 発話区間が重複する重複発話区間を検出し、

前記重複発話区間が所定の第 2 閾値未満の前記第 1 信号強度と前記第 2 信号強度以外の前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との相関係数を算出し、

前記相関係数に基づいて、前記第 1 音声と前記第 2 音声がか会話している状態か否かを判定する

ことを実行させることを特徴とする音声処理プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、例えば、音声処理装置、音声処理方法および音声処理プログラムに関する。

10

20

30

40

50

【背景技術】

【0002】

近年、例えば、企業内の従業員同士が円滑なコミュニケーションを実現できているか否かを把握する為に、録音機器（例えばマイクロフォン）を従業員毎に装着して、各従業員の音声を常時録音することが行われている。従業員の音声の常時録音によって、会議や雑談などの対面での会話音声（音声データと称しても良い）や、電話などの通信を介した会話音声など、様々なコミュニケーションにおける会話音声を取得することが出来る。常時録音された音声データから誰と誰がどれ位の時間に渡って会話したのかを分析することで、企業内の従業員同士が円滑なコミュニケーションを実現できているか否かを把握することが可能となる。

10

【先行技術文献】

【特許文献】

【0003】

【特許文献1】特開2010-266522号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

常時録音された会話音声から誰と誰がどれ位の時間会話したのかを分析する場合、複数の話者の音声個別に録音された音声データから、実際に会話が行われている話者に対応した音声データの組み合わせを判定する必要がある。この場合、評価者による主観評価により手作業で音声データの組み合わせの判定作業を実施する必要がある。また、組み合わせの判定作業を行う音声データが多い場合、音声データに対応する話者を特定することが評価者にとって困難な場合も想定され得る。一方、複数の話者の音声個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを判定することが可能となる音声処理装置は、実現されていない状況である。本発明においては、複数の話者の音声個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを判定することが可能となる音声処理装置を提供することを目的とする。

20

【課題を解決するための手段】

【0005】

本発明が開示する音声処理装置は、第1音声が含まれる第1入力信号と、第2音声が含まれる第2入力信号を取得する取得部と、第1入力信号の第1信号強度と、第2入力信号の第2信号強度を検出する検出部を備える。更に、当該音声処理装置は、第1信号強度の時系列と第2信号強度の時系列との相関係数を算出する算出部と、相関係数に基づいて、第1音声と第2音声が会話している状態か否かを判定する判定部を備える。

30

【0006】

なお、本発明の目的及び利点は、例えば、請求項におけるエレメント及び組み合わせにより実現され、かつ達成されるものである。また、上記の一般的な記述及び下記の詳細な記述の何れも、例示的かつ説明的なものであり、請求項のように、本発明を制限するものではないことを理解されたい。

【発明の効果】

40

【0007】

本明細書に開示される音声処理装置では、複数の話者の音声個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを判定することが可能となる。

【図面の簡単な説明】

【0008】

【図1】第1の実施形態による音声処理装置の機能ブロック図である。

【図2】音声処理装置による音声処理方法のフローチャートである。

【図3】一つの実施形態による検出部の機能ブロック図である。

【図4】実施例1に係る検出部3の検出処理のフローチャートである。

50

【図5】(a)は、検出部3による第1信号強度の第1検出結果を示す図である。(b)は、検出部3による第2信号強度と第2信号強度の第2検出結果を示す図である。

【図6】(a)は、検出部3による第1信号強度の参考検出結果を示す図である。(b)は、検出部3による第2信号強度と第2信号強度の参考検出結果を示す図である。

【図7】発話時間の比率の分布図である。

【図8】第1音声と第2音声が発話している状態における第1信号強度の時系列と第2信号強度の時系列との相関関係を示す図である。

【図9】第1信号強度の時系列と第2信号強度の時系列との相関係数の分布図である。

【図10】実施例1に係る算出部4の算出処理と判定部5の判定処理のフローチャートである。

10

【図11】実施例1と上述の比較例による複数の話者の音声が発話された音声データから、会話をしている音声データの組み合わせの判定性能を示す図である。

【図12】相関係数算出区間の概念図である。

【図13】相関係数算出区間に含まれる話者交替回数と、会話している状態または会話していない状態における当該相関係数算出区間に基づいて算出した相関係数の関係を示す図である。

【図14】信号強度の組み合わせに対する相関係数のテーブルを示す図である。

【図15】実施例4に係る検出部3の検出処理と算出部4の算出処理のフローチャートである。

【図16】実施例5に係る音声処理装置1の音声処理のフローチャートである。

20

【図17】(a)は、検出部3による第2信号強度の第2検出結果を示す図である。(b)は、検出部3による第3信号強度の第3検出結果を示す図である。(c)は、検出部3による第4信号強度の第4検出結果を示す図である。

【図18】一つの実施形態による音声処理装置として機能するコンピュータのハードウェア構成図である。

【発明を実施するための形態】

【0009】

以下に、一つの実施形態による音声処理装置、音声処理方法及び音声処理プログラムの実施例を図面に基づいて詳細に説明する。なお、当該実施例は、開示の技術を限定するものではない。

30

【0010】

(実施例1)

図1は、第1の実施形態による音声処理装置1の機能ブロック図である。音声処理装置1は、取得部2、検出部3、算出部4、判定部5を有する。図2は、音声処理装置1の音声処理のフローチャートである。実施例1においては、図2に示す音声処理装置1による音声処理のフローを、図1に示す音声処理装置1の機能ブロック図の各機能の説明に対応付けて説明する。

【0011】

図1の取得部2は、例えば、ワイヤードロジックによるハードウェア回路である。また、取得部2は、音声処理装置1で実行され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っても良い。取得部2は、入力音声の一例となる第1ユーザの第1音声と第2ユーザの第2音声を、例えば、外部装置を介して取得する。なお、当該処理は、図2に示すフローチャートのステップS201に対応する。取得部2は、第1音声が含まれる第1入力信号と第2音声が含まれる第2入力信号の取得方法として、例えば、特開2008-134557号公報に記載の方法を用いることが出来る。ここで、第1ユーザと第2ユーザは、図示しないマイクロフォン(上述の外部装置に相当)をそれぞれ装着しているものとする。第1音声は、例えば、第1ユーザの会話の相手となるユーザ(第2ユーザまたは第2ユーザ以外のユーザ)に対して発話する音声を示す第1ユーザの音声である。また、第2音声は、例えば、第2ユーザの会話の相手となるユーザ(第1ユーザまたは第1ユーザ以外のユーザ)に対して発話する音声を示す第2ユーザの音声である。また、第

40

50

1音声と第2音声は、例えば、日本語であるが、英語等の他の言語であっても良い。換言すると、実施例1における音声処理においては、言語依存は存在しない。なお、実施例1においては、説明の便宜上、取得部2は第1ユーザの第1音声が含まれる第1入力信号と第2ユーザの第2音声が含まれる第2入力信号を取得するものとして説明するが、取得部2は第1ユーザまたは第2ユーザ以外の第xユーザ(xは、例えば、音声処理装置1を使用するユーザ数)の第x音声が含まれる第x入力信号を取得しても良い。取得部2は取得した第1音声と第2音声を検出部3に出力する。

【0012】

検出部3は、例えば、ワイヤードロジックによるハードウェア回路である。また、検出部3は、音声処理装置1で実行され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っても良い。検出部3は、第1入力信号と第2入力信号を取得部2から受け取る。検出部3は、例えば、第1入力信号または第2入力信号に含まれる複数のフレームから、第1信号強度または第2信号強度を検出する。当該処理は、図2に示すフローチャートのステップS202に対応する。具体的には、検出部3は、第1入力信号に含まれる複数のフレーム毎の信号強度を示す第1信号強度を検出する。また、検出部3は、第2音声に含まれる複数のフレーム毎の信号強度を示す第2信号強度を検出する。なお、第1信号強度または第2信号強度は、例えば、第1音声または第2音声のパワーまたは信号対雑音比であれば良いが、当該パワーと信号対雑音比に限定されるものではない。

【0013】

ここで、検出部3による第1入力信号または第2入力信号のパワーまたは信号対雑音比の検出処理の詳細について説明する。なお、説明の便宜上、第1音声のパワーを第1パワー、第1音声の信号対雑音比を第1信号対雑音比と称し、第2音声のパワーを第2パワー、第2音声の信号対雑音比を第2信号対雑音比と称するものとする。なお、第2パワーの検出方法は、第1パワーの検出方法と同様の手法を用いることができ、第2信号対雑音比の検出方法は、第1信号対雑音比の検出方法と同様の手法を用いることができる。この為、実施例1においては、検出部3による第1パワーと第1信号対雑音比の検出処理の詳細について説明する。

【0014】

図3は、一つの実施形態による検出部3の機能ブロック図である。検出部3は、パワー算出部10、雑音推定部11、信号対雑音比算出部12を有する。なお、検出部3は、パワー算出部10、雑音推定部11、信号対雑音比算出部12を必ずしも有する必要はなく、各部が有する機能を、一つのまたは複数のワイヤードロジックによるハードウェア回路で実現させても良い。また、検出部3に含まれる各部が有する機能をワイヤードロジックによるハードウェア回路に代えて、音声処理装置1で実行され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実現させても良い。

【0015】

(検出部3によるパワー算出方法)

図3において、第1入力信号がパワー算出部10に入力される。なお、パワー算出部10は、図示しないバッファまたはキャッシュを有しても良い。パワー算出部10は、第1入力信号に含まれる各フレームの音量を算出し、当該音量を雑音推定部11と信号対雑音比算出部12へ出力する。なお、第1入力信号に含まれる各フレーム長は、160サンプル(8kHzサンプリングで20msに相当)である。第1音声信号の信号強度の一例となる第1パワー $P_1(t)$ は、次式の通り、第1信号 $s_1(t)$ のフレーム内の振幅二乗和に対して対数変換することで、算出することが出来る。

(数1)

$$P_1(t) = 10 \log \left(\sum_{i=t}^{t+L} s_1^2(i) \right)$$

上述の(数1)において、tはフレーム番号を示し、Lはフレーム長を示す。フレーム

10

20

30

40

50

長Lは、上述の通り、例えば160サンプル(8kHzサンプリングで20msに相当)である。

【0016】

(検出部3による信号対雑音比算出方法)

雑音推定部11は、各フレームの第1パワー $P_1(t)$ をパワー算出部10から受け取る。雑音推定部11は、各フレームにおける雑音を推定して、雑音推定結果を信号対雑音比算出部12へ出力する。ここで、雑音推定部11による各フレームの雑音推定は、例えば、以下の雑音推定方法を用いることが出来る。

【0017】

(雑音推定方法)

雑音推定部11は、第1入力信号のパワー $P_1(t)$ と、1フレーム過去の第1雑音電力 $N_1(t-1)$ の差に応じて第1雑音電力を次式に基づいて更新する。雑音推定部11は、第1音声信号のパワー $P_1(t)$ と、1フレーム過去の第1雑音電力 $N_1(t-1)$ の差が所定閾値(例えば、5dB)を下回る場合は雑音パワーを音声信号が雑音であると推定し、第1雑音電力 $N_1(t)$ を更新する。一方で、雑音推定部11は、第1入力信号のパワー $P_1(t)$ と、1フレーム過去の第1雑音電力 $N_1(t-1)$ の差が所定閾値以上の場合は第1雑音電力を次式に基づいて更新しない。

(数2)

$$N_1(t) = \begin{cases} M_1(t-1) * COF + P_1(t) * (1 - COF) & , (P_1(t) - N_1(t-1)) < TH_P \text{ の場合} \\ N_1(t-1) & , \text{ (上記以外)} \end{cases}$$

10

20

上述の(数2)において、 TH_P は雑音区間と判定するための判定閾値であり、例えば5dBであれば良い。また、 COF は忘却係数であり、例えば0.05であれば良い。雑音推定部11は雑音推定結果となる各フレームの第1雑音電力 $N_1(t)$ を受け取る。

【0018】

図3において、信号対雑音比算出部12は、パワー算出部10から各フレームの第1パワー $P_1(t)$ を受け取り、雑音推定部11から雑音推定結果となる各フレームの第1雑音電力 $N_1(t)$ を受け取る。なお、信号対雑音比算出部12は、図示しないキャッシュまたはメモリを有しており、過去Lフレーム分の第1パワー $P_1(t)$ 、第1雑音電力 $N_1(t)$ を保持することができる。信号対雑音比算出部12は、次式を用いて、第1信号対雑音比 $SNR_1(t)$ を算出する。

(数3)

$$SNR_1(t) = P_1(t) - N_1(t)$$

30

【0019】

図4は、実施例1に係る検出部3の検出処理のフローチャートである。検出部3は、第1入力信号と第2入力信号を取得部2から受け取る(ステップS401)。次に、検出部3は、第1入力信号の第1パワーと第2入力信号の第2パワーを上述の(数1)に基づいて検出する(ステップS402)。次に、第1音声の第1雑音電力と第2音声の第2雑音電力を上述の(数2)に基づいて検出する(ステップS403)。そして、検出部3は、第1パワーと第2パワーを第1信号強度と第2信号強度として選択するか否かを決定する(ステップS404)。例えば、第1雑音電力と第2雑音電力がそれぞれ所定の閾値を下回る場合に、第1パワーと第2パワーを第1信号強度と第2信号強度として選択する。検出部3は、第1パワーと第2パワーを第1信号強度と第2信号強度として選択することを決定した場合(ステップS404 - Yes)は、第1パワーと第2パワーを第1信号強度と第2信号強度として算出部4に出力する(ステップS407)ことで、図4のフローチャートに示す検出処理を終了する。一方、検出部3は、第1パワーと第2パワーを第1信号強度と第2信号強度として選択することを決定しない場合(ステップS404 - No)は、検出部3は、第1音声の第1信号対雑音比と第2音声の第2信号対雑音比を上述の(

40

50

数 3) に基づいて検出する (ステップ S 4 0 5) 。次に、検出部 3 は、第 1 信号対雑音比と第 2 信号対雑音比を、第 1 信号強度または第 2 信号強度として算出部 4 に出力する (ステップ S 4 0 6) ことで、図 4 のフローチャートに示す検出処理を終了する。

【 0 0 2 0 】

図 5 (a) は、検出部 3 による第 1 信号強度の第 1 検出結果を示す図である。図 5 (a) の横軸は時間を示し、縦軸は第 1 信号強度の一例となる第 1 パワーを示している。図 5 (b) は、検出部 3 による第 2 信号強度と第 2 信号強度の第 2 検出結果を示す図である。図 5 (b) の横軸は時間を示し、縦軸は第 2 信号強度の一例となる第 2 パワーを示している。なお、図 5 (a) と図 5 (b) の縦軸のパワーと横軸の時間は同一スケールである。検出部 3 は検出した第 1 信号強度と第 2 信号強度を算出部 4 に出力する。

10

【 0 0 2 1 】

図 1 において、算出部 4 は、例えば、ワイヤードロジックによるハードウェア回路である。また、算出部 4 は、音声処理装置 1 で実行され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っても良い。算出部 4 は、第 1 信号強度と第 2 信号強度を検出部 3 から受け取る。算出部 4 は、第 1 信号強度と第 2 信号強度の時系列に対する相関係数を算出する。なお、当該処理は、図 2 に示すフローチャートのステップ S 2 0 3 に対応する。具体的には、算出部 4 は、次式に基づいて、第 1 信号強度と第 2 信号強度の時系列に対する相関係数 $corr.$ を算出する。

(数 4)

$$corr. = \frac{\sum_{i=Ts}^{Te} (P1'(i) - \overline{P1'(i)}) (P2(i) - \overline{P2(i)})}{\sqrt{\sum_{i=Ts}^{Te} (P1'(i) - \overline{P1'(i)})^2} \sqrt{\sum_{i=Ts}^{Te} (P2(i) - \overline{P2(i)})^2}}$$

20

なお、上述の (数 4) において、 Ts は相関係数を算出する始点時刻を示し、 Te は相関係数を算出する終点時刻を示す。なお、 $Te - Ts$ は相関係数を算出する時間 (時系列と称しても良い) となるが、任意の長さの時間を適用することが可能であり、例えば、60 秒であれば良い。算出部 4 は算出した第 1 信号強度と第 2 信号強度の時系列に対する相関係数を判定部 5 に出力する。

【 0 0 2 2 】

ここで、算出部 4 が、第 1 信号強度と第 2 信号強度の時系列に対する相関係数を算出する技術的意義について説明する。まず、初めに本発明者らが検討した比較例について説明する。なお、当該比較例は公知技術ではなく、本発明者らにより新たに検証された事項であることを付言する。

30

【 0 0 2 3 】

(比較例)

複数の話者の音声データが個別に録音された音声データに対し、会話が行われている話者に対応した音声データの組み合わせを、2 つの音声データの「発話時間の比率」に基づいて判定する方法を比較例として考える。図 6 (a) は、検出部 3 による第 1 信号強度の参考検出結果を示す図である。図 6 (a) の横軸は時間を示し、縦軸は第 1 信号強度の一例となる第 1 パワーを示している。図 6 (b) は、検出部 3 による第 2 信号強度と第 2 信号強度の参考検出結果を示す図である。図 6 (b) の横軸は時間を示し、縦軸は第 2 信号強度の一例となる第 2 パワーを示している。なお、図 6 (a) と図 6 (b) の縦軸のパワーと横軸の時間は同一スケールである。ここで、任意の閾値 (例えば、10 dB) 以上の信号強度を満たす時間を発話時間と定義し、図 6 (a) と図 6 (b) の場合における発話時間の比率 (第 1 音声の発話時間 / 第 2 音声の発話時間) を算出すると、発話時間の比率は 1 . 1 となる。また、図 5 (a) と図 5 (b) の場合における発話時間の比率を、図 6 (a) と図 6 (b) の場合と同様に算出すると発話時間の比率は 1 . 1 となる。

40

【 0 0 2 4 】

ここで、自然な会話においては、話者同士が交互に発話を行うことでコミュニケーショ

50

ンを図る為、一方の話者の発話量が多い時間帯は、他方の話者の発話量が少ない（相手の発話に耳を傾ける）ことが想定される。ここで、図5（a）、図5（b）と図6（a）、図6（b）の組み合わせを客観的に比較した場合、一方の話者の発話量が多い時間帯に、他方の話者の発話量が少ないのは、図5（a）、図5（b）の組み合わせであることが理解できる。図6（a）、図6（b）の組み合わせは、一方の話者の発話量が多い時間帯に、他方の話者の発話量も多くなっている為、一般的には第1ユーザと第2ユーザは、それぞれお互いに別のユーザと会話していることが想定される。

【0025】

図7は、発話時間の比率の分布図である。図7においては、20話者の音声データの組み合わせとなる190組において、実際に会話をしている状態の10組の音声データの組み合わせと、会話を行っていない状態の180組の音声データの組み合わせの発話時間の比率の分布図を示している。図7の発話時間の比率の分布図から理解できる通り、会話をしている音声データの組み合わせの発話時間の比率と、会話を行っていない音声データの組み合わせの発話時間の比率は同等程度であることが確認された。この為、複数の話者の音声データが個別に録音された音声データから会話が行われている話者に対応した音声データの組み合わせを、比較例に示した2つの音声データの「発話時間の比率」に基づいて判定することは難しいことが、本発明者らの検証により明らかになった。

10

【0026】

一方、自然な会話においては、話者同士が交互に発話を行うことでコミュニケーションを図る為、図5（a）、図5（b）に示される通り、一方の話者の発話量が多い時間帯は、他方の話者の発話量が少ないことが想定される。ここで、図5（a）、図5（b）における、第1信号強度と第2信号強度の時系列に対する相関関係に着目すると、一方の変数（例えば、第1信号強度）が増大するほど、他方の変数（例えば、第2信号強度）が減少する「負の相関」の関係性を有していることが理解できる。

20

【0027】

この為、本実施例1の算出部4が、当該負の相関の強さを示す相関係数を算出することで、後述する判定部5が、当該相関係数に基づいて、第1音声と第2音声がか会話している状態か否かを判定することができる。図8は、第1音声と第2音声がか会話している状態における第1信号強度の時系列と第2信号強度の時系列との相関関係を示す図である。図8に示す通り、第1音声と第2音声がか会話している状態においては、相関関係は負の相関（負の相関係数）を有することが理解できる。

30

【0028】

図9は、第1信号強度の時系列と第2信号強度の時系列との相関係数の分布図である。図9においては、第1ユーザと第2ユーザの音声データを11組（計31分）用い、会話している区間と、会話していない区間（別の話者と会話している区間）を評価者が主観評価によって判定した上で、第1信号強度と第2信号強度の時系列に対する相関係数をそれぞれ算出した。図9から理解できる通り、会話している区間と会話していない区間で相関係数の範囲が明確に異なっていることが確認された。図9の実験結果より、任意の閾値である第1閾値を例えば、 -0.4 （一般的に相関の強さが中程度の値）と規定し、当該第1閾値と第1信号強度と第2信号強度の時系列に対する相関係数を用いることで、第1音声と第2音声がか会話している状態か否（会話していない状態）かを判定することができる。

40

【0029】

図1において、判定部5は、例えば、ワイヤードロジックによるハードウェア回路である。また、判定部5は、音声処理装置1で実行され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っても良い。判定部5は、算出部4が算出した第1信号強度の時系列と第2信号強度の時系列との相関係数を算出部4から受け取る。判定部5は、相関係数に基づいて第1音声と第2音声がか会話している状態か否かを判定する。なお、当該処理は、図2に示すフローチャートのステップS204に対応する。判定部5は、相関係数が負であり、相関係数が上述の第1閾値未満（例えば、 $-0.4 > \text{相関係数} \geq -1.0$

50

) の場合に第 1 音声と第 2 音声が会話している状態であると判定する。換言すると、判定部 5 は、相関係数が第 1 閾値以上 (例えば、 $-0.4 \leq \text{相関係数} \leq +1.0$) の場合に、第 1 音声と第 2 音声がか話していない状態であると判定しても良い。判定部 5 は、判定結果を任意の外部装置 (例えば、図示しないコンピュータ) に出力する。

【0030】

判定部 5 は、必要に応じて、第 1 音声と第 2 音声がか話している状態と判定した相関係数の算出に用いられた、上述の (数 4) における $T_e - T_s$ の区間長を算出部 4 から受け取り、当該区間長を、第 1 ユーザと第 2 ユーザがか話した区間 (時間) として外部装置に出力することもできる。または、判定部 5 は、必要に応じて、第 1 音声と第 2 音声がか話している状態と判定した相関係数の算出に用いられた、上述の (数 4) における始点時刻の T_s と、終点時刻の T_e を算出部 4 から受け取り、当該 T_s を第 1 ユーザと第 2 ユーザの会話の始点時刻、 T_e を第 1 ユーザと第 2 ユーザの会話の終点時刻として外部装置に出力することもできる。

10

【0031】

図 10 は、実施例 1 に係る算出部 4 の算出処理と判定部 5 の判定処理のフローチャートである。算出部 4 は、第 1 信号強度と第 2 信号強度を検出部 3 から受け取る (ステップ S1001)。次に、算出部 4 は、第 1 信号強度と第 2 信号強度の時系列に対する相関係数を (数 4) に基づいて算出する (ステップ S1002)。算出部 4 は、算出した相関係数を判定部 5 に出力する (ステップ S1003) ことで、図 10 のフローチャートに示す算出部 4 の算出処理を終了する。

20

【0032】

判定部 5 は、算出部 4 が算出した相関係数を算出部 4 から受け取る (ステップ S1004)。判定部 5 は、相関係数が上述の第 1 閾値未満か否かを判定する (ステップ S1005)。相関係数が上述の第 1 閾値未満の場合 (ステップ S1005 - Yes)、判定部 5 は、第 1 音声と第 2 音声がか話している状態として判定し (ステップ S1006)、その判定結果を任意の外部装置に出力する (ステップ S1008) ことで、図 10 のフローチャートに示す判定部 5 の判定処理を終了する。一方、相関係数が上述の第 1 閾値以上の場合 (ステップ S1005 - No)、判定部 5 は、第 1 音声と第 2 音声がか話していない状態として判定し (ステップ S1007)、その判定結果を任意の外部装置に出力する (ステップ S1008) ことで、図 10 のフローチャートに示す判定部 5 の判定処理を終了する。

30

【0033】

図 11 は、実施例 1 と上述の比較例による複数の話者の音声がか個別に録音された音声データから、会話をしている音声データの組み合わせの判定性能を示す図である。図 11 は、20 話者の音声データの組み合わせ計 190 組 (会話している音声データの組み合わせ 10 組、会話していない音声データの組み合わせ 180 組の合計) において、会話している音声データの検出率 (会話している音声を検出した割合) と、当該検出した会話している音声データの正解率 (検出した音声が会話している状態である割合) を示している。なお、比較例においては、検出率が 100% となるように発話時間比が 1.2 以下の音声データの組み合わせを会話しているデータとして判定した。図 11 から理解できる通り、実施例 1 によれば、比較例に対して正解率が大幅に向上していることが確認された。実施例 1 における音声処理装置 1 に依れば、会話において一方が発話している場合、他方は発話しない特徴を利用し、異なるユーザ同士の音声信号強度の時系列に対する相関係数が負であり、かつ相関係数が上述の第 1 閾値未満の場合にユーザ同士がか話していると判定することが可能となる。換言すれば、実施例 1 における音声処理装置 1 は、複数の話者の音声がか個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを、高い精度で判定することが可能となる。

40

【0034】

(実施例 2)

図 1 の検出部 3 は、実施例 1 の検出処理に加えて、第 1 信号強度と第 2 信号強度の時系

50

列に対する大小関係を検出し、当該大小関係が反転する反転回数が第4閾値以上（例えば、第4閾値 = 6回）となる相関係数算出区間を検出して良い。また、図1の算出部4は、実施例1の算出処理に加えて、実施例2に係る検出部3が検出する相関係数算出区間における相関係数を算出して良い。以下、実施例2に係る音声処理装置1の音声処理の詳細について説明する。

【0035】

検出部3は、第1信号強度と第2信号強度の時系列に対する大小関係を検出する。具体的には、検出部3は、第1信号強度と第2信号強度（例えば、第1パワーと第2パワー）の差が所定の第5閾値（例えば、パワー差 = +20dB）以上となる第1状態と、第1信号強度と第2信号強度（例えば、第1パワーと第2パワー）の差が所定の第6閾値（例えば、パワー差 = -20dB）以下となる第2状態を検出する。検出部3は、第1状態から第2状態へ遷移した場合、または、第2状態から第1状態に遷移した場合を話者交代点（第1信号強度と第2信号強度の時系列に対する大小関係が反転する点と称しても良い）として検出する。検出部3が、話者交代点が所定の第4閾値（例えば、第4閾値 = 6回）となる相関係数算出区間を検出する。図12は、相関係数算出区間の概念図である。図12の横軸は第1音声または第2音声の時系列であり、縦軸は、第1信号強度と第2信号強度の差の一例となるパワー差である。なお、図12のパワー差は、図5(a)と図5(b)のパワー差を示している。検出部3は、話者交代点が第4閾値以上を満たす区間を相関係数算出区間として検出する。検出部3は、検出した相関係数算出区間を算出部4に出力する。

10

20

【0036】

算出部4は、検出部3が検出した相関係数算出区間を検出部3から受け取る。算出部4は、相関係数算出区間における第1信号強度と第2信号強度の時系列に対する相関係数を実施例1と同様の手法を用いて算出する。算出部4は、上述の(数4)のTsを相関係数算出区間の始点時刻とし、Teを相関係数算出区間の終点時刻とし、Te - Tsで表現される相関係数算出区間における第1信号強度と第2信号強度の時系列に対する相関係数を算出すれば良い。なお、算出部4は、検出部3から複数の相関係数算出区間を受け取っている場合は、複数の相関係数算出区間毎に相関係数を算出すれば良い。算出部4は、算出した相関係数を判定部5に出力し、判定部5は実施例1と同様の方法で判定処理を行えば良い。なお、判定部5は、複数の相関係数を算出部4から受け取っている場合は、実施例1と同様の方法で相関係数毎に判定処理を行えば良い。

30

【0037】

ここで、実施例2における技術的意義について説明する。図13は、相関係数算出区間に含まれる話者交替回数と、会話している状態または会話していない状態における当該相関係数算出区間に基づいて算出した相関係数の関係を示す図である。図13から理解できる通り、話者交代回数が多い区間（例えば、上述の第4閾値となる6回）の場合、会話している音声データの組み合わせの相関係数と、会話していない音声データの組み合わせの差（有意差）が現れる。この為、実施例2における音声処理装置1においては、話者交代回数が所定の第4閾値以上となる相関係数算出区間で相関係数を算出することで、高い精度で複数の話者の音声個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを、高い精度で判定することが可能となる。なお、実施例2における音声処理装置1は、実施例1に記載した音声処理を任意に組み合わせることができる。

40

【0038】

（実施例3）

図1の算出部4は、実施例1または実施例2の算出処理に加えて、第1信号強度の第1位相または、第2信号強度の第2位相を所定の範囲で変化させて複数の相関係数を算出しても良い。また、図1の判定部5は、算出部4が第1位相または第2位相を所定の範囲で変化させて算出した複数の相関係数の中で、最小値を満たす相関係数に基づいて、第1音声と第2音声が会話している状態か否かを判定しても良い。以下、実施例3に係る音声処理装置1の音声処理の詳細について説明する。

50

【 0 0 3 9 】

算出部 4 は、第 1 信号強度の第 1 位相または、第 2 信号強度の第 2 位相を所定の範囲で変化させて複数の相関係数を算出する。実施例 3 においては、第 1 位相を所定の範囲で変化させた場合について説明するが、第 2 位相を所定の範囲で変化させる場合も第 1 位相の場合と同様に算出することが出来る為、第 2 位相に関する詳細な説明は省略する。算出部 4 は、次式に基づいて、第 1 信号強度（例えば、第 1 パワー P 1）の位相を変化させた場合の第 1 信号強度 P 1（t）と第 2 信号強度 P 2（t）の複数の相関係数から、最小値の相関係数を満たす位相 d m i n を次式に基づいて算出する。

(数 5)

$$dmin = \operatorname{argmin}_d \left\{ \frac{\sum_{i=Ts}^{Te} (P1(i) + d - \overline{P1(i)}) (P2(i) - \overline{P2(i)})}{\sqrt{\sum_{i=Ts}^{Te} (P1(i) + d - \overline{P1(i)})^2} \sqrt{\sum_{i=Ts}^{Te} (P2(i) - \overline{P2(i)})^2}} \right\} \mid d = -Dmax, \dots, Dmax$$

10

$$\overline{P1(i)} = \frac{1}{Te - Ts + 1} \sum_{i=Ts}^{Te} P1(i)$$

上述の（数 5）において d は位相の変更量（サンプル）を示し、D m a x は位相変更の最大値を示す。D m a x は、例えば 8 0 0 0 0 サンプル（10 秒に相当）とすれば良い。算出部 4 は、第 1 強度の位相を所定の範囲となる d m i n 変化（シフトと称しても良い）させた信号 P 1' t を次式に基づいて算出する。

20

(数 6)

$$P1'(t) = P1(t + dmin)$$

判定部 5 は、複数の相関係数のうち最小値の満たす相関係数に基づいて、第 1 音声と第 2 音声がかかっている状態か否かを実施例 1 と同様の判定方法を用いて判定する。

【 0 0 4 0 】

実施例 3 における音声処理装置 1 によれば、例えば、取得部 2 が取得する第 1 入力信号と第 2 入力信号のそれぞれにおいて、録音のタイミングが同期していない場合（例えば、第 1 ユーザと第 2 ユーザがそれぞれ装着しているマイクロフォンの音声処理の内部処理の差異や、設定時刻のずれ等に起因して録音のタイミングが同期していない場合）においても、相関係数を算出する時点において、位相を調整することによって、会話している音声データの組み合わせの判定精度を向上させることが可能となる。なお、実施例 3 における音声処理装置 1 は、実施例 1 または実施例 2 に記載した音声処理を任意に組み合わせることができる。

30

【 0 0 4 1 】

(実施例 4)

図 1 の検出部 3 は、実施例 1 ないし実施例 3 の検出処理に加え、第 1 入力信号に含まれる第 1 発話区間を第 1 信号強度に基づいて検出し、第 2 入力信号に含まれる第 2 発話区間を第 2 信号強度に基づいて検出し、第 1 発話区間と第 2 発話区間が重複する重複発話区間を検出しても良い。また、図 1 の算出部 4 は、実施例 1 ないし実施例 3 の算出処理に加え、当該重複発話区間が第 2 閾値未満の第 1 信号強度と第 2 信号強度を、相関係数の算出に用いない算出処理を行っても良い。また、検出部 3 は、第 1 入力信号に含まれる第 1 無音区間を第 1 信号強度に基づいて検出し、第 2 入力信号に含まれる第 2 無音区間を第 2 信号強度に基づいて検出し、第 1 無音区間と第 2 無音区間が重複する重複無音区間を検出しても良い。算出部 4 は、当該重複無音区間が第 3 閾値未満の第 1 信号強度と第 2 信号強度を、相関係数の算出に用いない算出処理を行っても良い。以下、実施例 4 に係る音声処理装置 1 の音声処理の詳細について説明する。

40

【 0 0 4 2 】

50

検出部 3 は、例えば、第 1 信号のフレーム単位で、第 1 信号強度が所定の閾値（例えば、10 dB（但し、第 1 信号強度として第 1 パワーを用いる場合））以上か否かを判定する。検出部 3 は、第 1 信号強度が当該閾値以上を満たすフレームを第 1 発話区間として判定する。また、検出部 3 は、第 1 信号強度が当該閾値未満を満たすフレームを第 1 無音区間として判定する。検出部 3 は、次式に基づいて、フレーム単位で当該フレームが第 1 発話区間か第 1 無音区間（非発話区間と称しても良い）を判定し、判定結果 $v_1(t)$ を検出する。

（数 7）

$$v_1(t) = 1 \quad (\text{第 1 発話区間})$$

$$v_1(t) = 0 \quad (\text{第 1 無音区間})$$

10

上述の（数 7）において、 t はフレーム番号を示す。なお、1 フレームの長さは、例えば、20 ms である。また、上述の（数 7）においては、 t フレーム目の第 1 音声が発話区間と判定された場合は $v_1(t) = 1$ が代入され、 t フレーム目の第 1 音声が無音区間と判定された場合は $v_1(t) = 0$ と代入されることを意味する。検出部 3 は、 $v_1(t) = 1$ を連続して満たすフレーム区間を第 1 発話区間として出力する。なお、検出部 3 は、第 2 音声に関する判定結果 $v_2(t)$ を $v_1(t)$ と同様の手法を用いて算出し、第 2 発話区間または第 2 無音区間を検出する。

20

【0043】

次に、検出部 3 は、第 1 発話区間と第 2 発話区間が重複する重複発話区間を検出する。当該重複発話区間は、例えば、ある任意の同一の時刻において、第 1 ユーザと第 2 ユーザが互いに発話している区間と定義することが出来る。なお、検出部 3 は、具体的には、次式に基づいて重複発話区間 $TO(t)$ を検出することが出来る。

（数 8）

$$\text{if } \{v_1(t) = 0\} \vee \{v_2(t) = 0\} \quad TO(t) = 0$$

$$\text{else} \quad TO(t) = TO(t-1) + 1$$

30

上述の（数 8）は、第 1 ユーザの第 1 音声と第 2 ユーザの第 2 音声の何れかが無音区間と判定されるフレームに対しては重複時間を 0（重複区間の発現無し）と規定し、第 1 ユーザの第 1 音声と第 2 ユーザの第 2 音声の双方が発話区間と判定されるフレームに対しては直前のフレームまでの重複時間に 1 フレーム加算することで、重複が連続するフレーム数（重複区間）を算出することを意味する。検出部 4 は、規定した重複発話区間を算出部 5 に出力する。なお、重複発話区間には、区間の長さの情報が含まれているものとする。重複発話区間の長さ LO は、例えば、次式に基づいて算出することができる。

（数 9）

$$LO = TO_e(i) - TO_s(i)$$

40

なお、上述の（数 9）において、 $TO_s(i)$ は、重複発話区間の始点（開始フレーム）であり、 $TO_e(i)$ は、重複発話区間の終点（終了フレーム）である。

【0044】

次に、検出部 3 は、第 1 無音区間と第 2 無音区間が重複する重複無音区間を検出する。当該重複無音区間は、例えば、ある任意の同一の時刻において、第 1 ユーザと第 2 ユーザが互いに発話していない区間と定義することが出来る。なお、検出部 3 は、具体的には、次式に基づいて重複無音区間 $TE(t)$ を検出することが出来る。

（数 10）

50

検出した重複発話区間と重複無音区間を算出部4に出力することで、図15のフローチャートに示す検出処理を終了する。

【0048】

算出部4は、重複発話区間と重複無音区間を検出部3から受け取り、重複発話区間が上述の第2閾値未満、かつ重複無音区間が上述の第3閾値未満か否かを判定する(ステップS1505)。重複発話区間が上述の第2閾値未満、かつ重複無音区間が上述の第3閾値未満の場合(ステップS1505 - Yes)、算出部4は、重複発話区間と重複無音区間を用いずに相関係数を上述の(数4)に基づいて算出し(ステップS1509)、算出した相関係数を判定部5に出力する(ステップS1512)ことで、図15のフローチャートに示す算出処理を終了する。

10

【0049】

重複発話区間が上述の第2閾値未満、かつ重複無音区間が上述の第3閾値未満ではない場合(ステップS1505 - No)、算出部4は、重複発話区間が第2閾値未満、かつ重複無音区間が第3閾値以上か否かを判定する(ステップS1506)。重複発話区間が第2閾値未満、かつ重複無音区間が第3閾値以上の場合(ステップS1506 - Yes)、算出部4は、重複発話区間を用いずに相関係数を上述の(数4)に基づいて算出し(ステップS1510)、算出した相関係数を判定部5に出力する(ステップS1512)ことで、図15のフローチャートに示す算出処理を終了する。

【0050】

重複発話区間が上述の第2閾値未満、かつ重複無音区間が上述の第3閾値以上ではない場合(ステップS1506 - No)、算出部4は、重複発話区間が第2閾値以上、かつ重複無音区間が第3閾値未満か否かを判定する(ステップS1507)。重複発話区間が第2閾値以上、かつ重複無音区間が第3閾値未満の場合(ステップS1507 - Yes)、算出部4は、重複無音区間を用いずに相関係数を上述の(数4)に基づいて算出し(ステップS1511)、算出した相関係数を判定部5に出力する(ステップS1512)ことで、図15のフローチャートに示す算出処理を終了する。

20

【0051】

重複発話区間が第2閾値以上、かつ重複無音区間が第3閾値未満ではない場合(ステップS1507 - No)、算出部4は、実施例1と同様の方法で、相関係数を上述の(数4)に基づいて算出し(ステップS1508)、算出した相関係数を判定部5に出力する(ステップS1512)ことで、図15のフローチャートに示す算出処理を終了する。

30

【0052】

実施例4においては、説明の便宜上、重複発話区間と重複無音区間に対する双方の処理について説明したが、重複発話区間または重複無音区間の何れか一方のみを処理対象としても良いし、双方を処理対象としても良い。実施例4に係る音声処理装置においては、複数の話者の音声が入力された音声データから、会話が行われている話者に対応した音声データの組み合わせを、重複発話区間または重複無音区間を考慮することで、高い精度で判定することが可能となる。なお、実施例4における音声処理装置1は、実施例1ないし実施例3に記載した音声処理を任意に組み合わせることができる。

【0053】

(実施例5)

図1の取得部2は、第1ユーザの第1音声が含まれる第1入力信号または第2ユーザの第2音声が含まれる第1入力信号のみならず、複数のユーザ(第xユーザと称しても良い)の各音声が含まれる各入力信号(第x入力音声と称しても良い)を取得することが出来る。実施例5においては、複数のユーザの各音声の組み合わせにおいて、実際に会話をしている2つの音声データの組み合わせを判定する方法について説明する。なお、実施例5は、説明の便宜上、第1ユーザの第1音声ないし第4ユーザの第4音声を取得部2が外部装置を介して取得するものとして説明する。ここで、第1ユーザないし第4ユーザは、図示しないマイクロフォン(上述の外部装置に相当)をそれぞれ装着しているものとする。取得部2は取得した第1音声ないし第4音声を検出部3に出力する。

40

50

【 0 0 5 4 】

図 1 の検出部 3 は、第 1 入力音声ないし第 4 入力音声を取得部 3 から受け取る。検出部 3 は、第 1 信号ないし第 4 音声の各信号強度（第 1 信号強度、第 2 信号強度、第 3 信号強度、第 4 信号強度）を、例えば、実施例 1 に記載した検出方法を用いて検出する。検出部 3 は、検出した第 1 信号強度ないし第 4 信号強度を算出部 4 に出力する。

【 0 0 5 5 】

図 1 の算出部 4 は、各信号強度、換言すると第 1 信号強度ないし第 4 信号強度を検出部 3 から受け取る。算出部 4 は、各信号強度から 2 つの信号強度を組み合わせた場合の時系列に対する各相関係数を算出する。具体的には、算出部 4 は、第 1 信号強度ないし第 4 信号強度から 2 つ信号強度を組み合わせる。図 1 4 は、信号強度の組み合わせに対する相関係数の 10
 テーブルを示す図である。算出部 4 は、例えば、算出部 4 が有する図示しないバッファまたはメモリに図 1 4 に示す相関係数のテーブルを格納し、音声データの組み合わせ ID に対応付けて第 1 信号強度ないし第 4 信号強度から 2 つ信号強度を組み合わせる。算出部 4 は、図 1 4 の組み合わせ ID に対応付けられている、2 つの信号強度を組み合わせた場合の、当該 2 つの信号強度の時系列に対する各相関係数を、例えば、実施例 1 で記載した算出方法を用いて算出する。算出部 4 は、例えば、図 1 4 の相関係数のテーブルに、各相関係数を格納しても良い。算出部 4 は、算出した各相関係数を判定部 5 に出力する。

【 0 0 5 6 】

図 1 の判定部 5 は、第 1 信号強度ないし第 4 信号強度から 2 つの信号強度を組み合わせた場合の時系列に対する各相関係数を算出部 4 から受け取る。判定部 5 は、各相関係数のうち、相関係数が最小値を満たす 2 つの信号強度の組み合わせに基づいて、複数のユーザの中から会話しているユーザの組み合わせを判定する。例えば、図 1 4 の相関係数の組み合わせ 20
 テーブルを参照すると、組み合わせ ID 4 における第 2 信号強度と第 3 信号強度が最小値を満たす相関係数である為、判定部 5 は、第 2 ユーザと第 3 ユーザが会話している状態のユーザの組み合わせとして判定する。また、判定部 5 は、第 2 信号強度と第 3 信号強度が最小値となる場合において、第 2 信号強度と第 3 信号強度の時系列に対する相関係数が上述の第 1 閾値未満の場合に、組み合わせた 2 つの音声データが会話している状態として判定しても良い。

【 0 0 5 7 】

図 1 6 は、実施例 5 に係る音声処理装置 1 の音声処理のフローチャートである。取得部 2 は、第 1 音声ないし第 4 音声（複数の音声と称しても良い）を取得する（ステップ S 1 6 0 1）。次に、検出部 3 は、第 1 信号強度ないし第 4 信号強度を検出する（ステップ S 1 6 0 2）。次に、検出部 3 は、第 1 信号強度ないし第 4 信号強度から 2 つの信号強度を組み合わせる（ステップ S 1 6 0 3）。次に、算出部 4 は、組み合わせた 2 つの信号強度の時系列に対する相関係数を（数 4 に）基づいて算出する（ステップ S 1 6 0 4）。次に、算出部 4 は、組み合わせた全ての 2 つの信号強度の時系列に対する相関係数を算出したか否かを判定する（ステップ S 1 6 0 5）、組み合わせた全ての 2 つの信号強度の時系列に対する相関係数を算出していない場合（ステップ S 1 6 0 5 - N o）、算出部 4 は、ステップ S 1 6 0 4 と S 1 6 0 5 の処理を繰り返し実施する。組み合わせた全ての 2 つの信号強度の時系列に対する相関係数を算出している場合（ステップ S 1 6 0 5 - Y e s）、算出部 4 は 40
 、相関係数が最小値を満たす 2 つの信号強度を判定する（ステップ S 1 6 0 6）。

【 0 0 5 8 】

次に、判定部 5 は、最小値を満たす相関係数が上述の第 1 閾値未満か否かを判定する（ステップ S 1 6 0 7）。相関係数が上述の第 1 閾値未満の場合（ステップ S 1 6 0 7 - Y e s）、判定部 5 は、相関係数が最小値を満たす 2 つの信号強度が会話している状態として判定し（ステップ S 1 6 0 8）、その判定結果を任意の外部装置に出力する（ステップ S 1 6 1 0）ことで、音声処理装置 1 は、図 1 6 のフローチャートに示す音声処理を終了する。相関係数が上述の第 1 閾値以上の場合（ステップ S 1 6 0 7 - N o）、判定部 5 は、組み合わせた 2 つの信号強度の何れも会話していない状態として判定し、（ステップ S 1 6 0 9）、その判定結果を任意の外部装置に出力する（ステップ S 1 6 1 0）ことで、音 50

声処理装置 1 は、図 16 のフローチャートに示す音声処理を終了する。

【0059】

実施例 5 における音声処理装置 1 によれば、複数の話者の音声は個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを、相関係数が最小値を満たす 2 つの信号強度の組み合わせに基づいて判定することで、高い精度で判定することが可能となる。なお、実施例 5 における音声処理装置 1 は、実施例 1 ないし実施例 4 に記載した音声処理を任意に組み合わせることができる。

【0060】

(実施例 6)

実施例 6 における音声処理装置 1 においては、3 人以上で会話しているグループを特定することが可能である。例えば、職場等においては、会話している話者は 2 名とは限らず、3 名以上のグループで会話をしている場合も存在する。以下、3 人以上で会話しているグループを特定する音声処理装置 1 の音声処理の詳細について説明する。実施例 6 における取得部 2、検出部 3 の処理、ならびに算出部 4 の一部の処理は、実施例 5 と同様の処理である為、重複する処理に関する詳細な説明は省略する。

【0061】

図 1 の判定部 5 は、第 1 信号強度ないし第 4 信号強度から 2 つの信号強度を組み合わせた場合の時系列に対する各相関係数を算出部 4 から受け取る。判定部 5 は、実施例 5 と同様に、各相関係数のうち、相関係数が最小値を満たす 2 つの信号強度の組み合わせに基づいて、複数のユーザの中から会話している信号強度の組み合わせを判定する。例えば、図 14 の相関係数の組み合わせテーブルを参照すると、組み合わせ ID 4 における第 2 信号強度と第 3 信号強度が最小値を満たす相関係数の組み合わせになる。次に、算出部 4 は、最小値を満たす相関係数の算出に用いた 2 つの信号強度を加算した加算信号強度を算出する。具体的には、例えば、算出部 4 は第 2 信号強度 $P_2(t)$ と第 3 信号強度 $P_3(t)$ を加算した加算信号強度 $P_{ADD}(t)$ を次式 (数 12) または (数 13) の何れか一方、または双方の平均値を用いても良い) に基づいて算出する。

(数 12)

$$P_{ADD_n}(t) = \sum_i P_i(t) \mid i \in G_n$$

(数 13)

$$P_{ADD_n}(t) = \max_i P_i(t) \mid i \in G_n$$

上述の (数 12) と (数 13) において、 n は第 x ユーザに対応する。また、上述の (数 12) は、信号強度の合計値を示し、(数 13) は信号強度の最大値を示す。また、 G_n はグループを示す。例えば、第 1 信号と第 2 信号がグループに含まれる場合は、 $G_n = \{1, 2\}$ と表現され、第 2 信号と第 3 信号がグループに含まれる場合は、 $G_n = \{2, 3\}$ と表現される。

【0062】

次に、算出部 4 は、加算信号強度と、最小値を満たす相関係数の算出に用いた 2 つの信号強度以外 (例えば、第 2 信号強度または第 3 信号強度以外) 1 つの信号強度 (第 1 信号強度または第 4 信号強度) との時系列に対する参照相関係数を算出する。実施例 6 においては、説明の便宜上、加算信号強度と第 4 ユーザの第 4 信号強度の時系列に対する参照相関係数 $corr_n$ を次式に基づいて算出する。

(数 14)

10

20

30

40

50

$$corr_n = \frac{\sum_{i=Ts}^{Te} (P1(i) - \overline{P1(t)}) (P_ADDn(i) - \overline{P_ADD(t)})}{\sqrt{\sum_{i=Ts}^{Te} (P1(i) - \overline{P1(t)})^2} \sqrt{\sum_{i=Ts}^{Te} (P_ADD(i) - \overline{P_ADD(t)})^2}}$$

上述の(数14)において、実施例6においては、 $P1(i)$ は、第4信号強度であり、 $P_ADD(t)$ は、加算信号強度である。なお、算出部4は、最小値を満たす相関係数の算出に用いた2つの信号強度の重複無音区間が第5閾値以上(例えば、重複無音区間 = 10秒)の場合に、参照相関係数を算出して良い。算出部4は、2つの信号強度の組み合わせの中で最小値を満たす相関係数と、参照相関係数を判定部5に出力し、判定部5はそれらを受け取る。

10

【0063】

判定部5は、参照相関係数と2つの信号強度の組み合わせの中で最小値を満たす相関係数を算出部4から受け取る。判定部5は、参照相関係数が最小値を満たす相関係数未満の場合、または、参照相関係数が上述の第1閾値未満の場合に、参照相関係数の算出に用いた3つの信号強度に基づいて、会話している3人以上のユーザの組み合わせ(会話しているグループと称しても良い)を判定する。例えば、判定部5は、参照相関係数が最小値を満たす相関係数未満の場合、第2ユーザと第3ユーザと第4ユーザが3名のグループで会話しているものと判定する。また、判定部5は、参照相関係数が、最小値を満たす相関係数未満の場合、かつ、第1閾値未満の場合、第2ユーザと第3ユーザと第4ユーザが3名のグループで会話しているものと判定しても良い。更に、判定部5は、参照相関係数が、最小値を満たす相関係数以上の場合であっても、第1閾値未満の場合であれば、第2ユーザと第3ユーザと第4ユーザが3名のグループで会話しているものと判定しても良い。なお、例えば、上述の(数12)または(数13)において、グループ $G_n = \{2, 3\}$ であり(第2信号と第3信号がグループに含まれる場合)、第4音声を加える場合は、新たなグループとして、 $G_{n+1} = G_n \cup \{3\} = \{2, 3, 4\}$ が定義されれば良い。

20

【0064】

また、判定部5が、第2ユーザと第3ユーザと第4ユーザが3名のグループで会話しているか否かを判定した後、以下の処理を実施することで、第1ユーザを含めた4名のグループで会話しているか否かを判定することができる。まず、算出部4は、第2信号強度、第3信号強度、第4信号強度に基づいて算出した参照相関係数を、最小値を満たす相関係数に置換する。次に、算出部4は、第1信号強度、第2信号強度、ならびに第3信号強度の上述の(数12)または、(数13)に基づいて、加算信号強度を算出する。次に、次に、算出部4は、第1信号強度と、加算信号強度の時系列に対する参照相関係数を算出する。判定部5の処理は、実施例6に開示している同様の判定方法を用いれば良い。判定部5は、全ての信号強度に対して上述の処理を行い、会話をしている複数のユーザ(グループ)を判定することもできる。

30

【0065】

ここで、実施例6の技術的意義について説明する。図17(a)は、検出部3による第2信号強度の第2検出結果を示す図である。図17(a)の横軸は時間を示し、縦軸は第2信号強度の一例となる第2パワーを示している。図17(b)は、検出部3による第3信号強度の第3検出結果を示す図である。図17(b)の横軸は時間を示し、縦軸は第3信号強度の一例となる第3パワーを示している。図17(c)は、検出部3による第4信号強度の第4検出結果を示す図である。図17(c)の横軸は時間を示し、縦軸は第4信号強度の一例となる第4パワーを示している。なお、図17(a)ないし図17(c)の縦軸のパワーと横軸の時間は同一スケールである。

40

【0066】

2名の話者における自然な会話においては、話者同士が交互に発話を行うことでコミュ

50

ニケーションを図る為、実施例1で説明した通り、一方の話者の発話量が多い時間帯は、他方の話者の発話量が少ないことが想定される。ここで、3名の話者が存在し、3名の話者同士が交互に発話を行う場合、1名が発話している場合、残りの2名の発話量が少ないことが想定される。この場合において、実質的に3名の音声の信号強度の相関係数を算出することで、より高い精度で会話しているユーザの組み合わせを判定することが出来る。例えば、図17において、第2信号強度と第3信号強度が0の時間帯（重複無音区間に相当）が存在する。この時間帯に第1信号強度が0より大きければ、参照相関係数は、第2信号強度と第3信号強度の相関係数未満になる。また、第2信号強度と第3信号強度の重複無音区間が第5閾値以上（例えば、重複無音区間 = 10秒）の場合（換言すると、第1信号強度が0より大きいことが想定される場合）のみに、参照相関係数を算出する処理を実施しても良い。当該処理により算出部4の算出処理の負荷を軽減することができる。また、重複無音区間において、信号強度が大きいフレームを最も多く含む音声の信号強度を、優先的に参照相関係数を算出する信号強度として選択しても良い。

【0067】

実施例6における音声処理装置1によれば、複数の話者の音声は個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを、高い精度で判定することが可能となる。なお、実施例6における音声処理装置1は、実施例1ないし実施例5に記載した音声処理を任意に組み合わせることができる。

【0068】

（実施例7）

図1の検出部3は、実施例1ないし実施例6の検出処理に加え、第1入力信号に含まれる第2信号強度または、第2入力信号に含まれる第1信号強度を更に検出することができる。また、算出部4は、実施例1ないし実施例6の算出処理に加え、第1入力信号に含まれる第2信号強度と第2入力信号に含まれる第2信号強度、または、第2信号に含まれる第1信号強度と第1信号に含まれる第1信号強度の時系列に対する第2相関係数を更に算出することができる。更に、判定部5は、実施例1ないし実施例6の判定処理に加え、第2相関係数に基づいて、第1音声と第2音声が発話している状態が否かを判定することができる。以下、実施例7に係る音声処理装置1の音声処理の詳細について説明する。

【0069】

例えば、第1ユーザと第2ユーザが会議や雑談などの対面での会話している場合、双方の距離間隔は比較的近接している為、第1ユーザが装着しているマイクロフォンには、第1ユーザの第1音声のみならず、第2ユーザの第2音声も入力される場合もある。同様に、第2ユーザが装着しているマイクロフォンには、第2ユーザの第2音声のみならず、第1ユーザの第1音声も入力される場合もある。この為、実施例7においては、検出部3が取得部2から受ける第1信号には第2信号も含まれているものとし、第2信号には第1信号も含まれているものとして説明する。なお、説明の便宜上、第1信号に含まれる第2信号の区間、または、第2信号に含まれる第1信号の区間を回り込み区間と称するものとする。

【0070】

検出部3は、第1信号に含まれる第2信号を、第1信号から分離する。検出部3は、具体的には、第1音声の第1信号強度の一例となる第1パワー $P_1(t)$ と、第2音声の第2信号強度の一例となる第2パワー $P_2(t)$ について、所定の時間間隔（例えば、1秒）の範囲毎に相関係数 $corr_n$ を算出し、相関性が高いと判定される場合に、当該範囲を回り込み区間と判定することができる。

（数15）

$$corr_n = \frac{\sum_{i=Tn-1}^{Tn} (P1(i) - \overline{P1(i)}) (P2(i) - \overline{P2(i)})}{\sqrt{\sum_{i=Tn-1}^{Tn} (P1(i) - \overline{P1(i)})^2} \sqrt{\sum_{i=Tn-1}^{Tn} (P2(i) - \overline{P2(i)})^2}}$$

10

20

30

40

50

if $corr_n > TH_SNEAK$: 回り込み区間 ($\overline{P1(i)}$ と $\overline{P2(i)}$)を比較し、小さい方が回り込み区間)
 else : 回り込み区間でない

ここで、上述の(数15)において、 T_n は相関算出範囲のフレーム長(例えば、1フレーム = 20 msec)を示し、例えば、1 secの場合は $T_n = n * 50$ とする。または、回り込み区間の判定閾値となる TH_SNEAK は、例えば0.95とすれば良い。また、検出部3は、同様に上述の(数15)を用いて、第2信号に含まれる第1信号を、次式に基づいて第2信号から分離することができる。検出部3は、第1信号に含まれており、当該第1信号から分離した第2信号、または、第2信号に含まれており、当該第2信号から分離した第1信号を算出部4に出力する。なお、説明の便宜上、第1信号に含まれており、当該第1信号から分離した第2信号を「第2分離信号」と称するものとする。また、第2信号に含まれており、当該第2信号から分離した第1信号を「第1分離信号」と称するものとする。検出部3は、第1分離信号と第2分離信号を算出部4に出力する。
 【0071】

算出部4は、第1分離信号と第2分離信号を検出部3から受け取る。算出部4は、第1分離信号の信号強度と第1信号強度、または、第2分離信号の信号強度と第2信号強度の信号強度の時系列に対する第2相関係数を、上述の(数4)を用いて算出する。算出部4は、算出した第2相関係数を判定部5に出力する。ここで、算出部4が第2相関係数を算出する技術的意義について説明する。例えば、第1分離信号の信号強度と第1信号強度の時系列に対する第2相関係数を考えると、第1分離信号と第1信号は、入力されるマイクロフォンが異なる為、信号強度自体は異なるが、第1ユーザから発せられたものである。この為、第1信号がある程度の強度を有する場合、第1分離信号もある程度の強度を有することになる。この為、例えば、第1分離信号の信号強度と第1信号強度の時系列に対する第2相関係数は正の相関を有することになる。この為、算出部4が、正の相関の強さを示す第2相関係数を算出することで、後述する判定部5が、当該第2相関係数に基づいて、第1音声と第2音声が会話している状態か否かを判定することができる。

【0072】

判定部5は、第2相関係数を算出部4から受け取る。判定部5は、第2相関係数に基づいて、第1音声と第2音声が会話している状態か否かを判定することができる。例えば、判定部5は、第2相関係数が正であり、当該第2相関係数が所定の第6閾値(例えば第6閾値 = +0.4)以上の場合に、第1音声と第2音声が会話している状態と判定することができる。また、判定部5は、第2相関係数が第6閾値以上の場合かつ、相関係数が第1未満の場合に第1音声と第2音声が会話している状態と判定することができる。

【0073】

実施例7における音声処理装置1によれば、複数の話者の音声は個別に録音された音声データから、会話が行われている話者に対応した音声データの組み合わせを、第2相関係数に基づいて判定することで、高い精度で判定することが可能となる。なお、実施例6における音声処理装置1は、実施例1ないし実施例6に記載した音声処理を任意に組み合わせることができる。

【0074】

(実施例8)

図18は、一つの実施形態による音声処理装置1として機能するコンピュータのハードウェア構成図である。図18に示す通り、音声処理装置1は、コンピュータ100、およびコンピュータ100に接続する入出力装置(周辺機器)を含んで構成される。

【0075】

コンピュータ100は、プロセッサ101によって装置全体が制御されている。プロセッサ101には、バス109を介してRAM(Random Access Memory)102と複数の周辺機器が接続されている。なお、プロセッサ101は、マルチプロセッサであってもよい。また、プロセッサ101は、例えば、CPU、MPU(Micro

10

20

30

40

50

o Processing Unit)、DSP(Digital Signal Processor)、ASIC(Application Specific Integrated Circuit)、またはPLD(Programmable Logic Device)である。更に、プロセッサ101は、CPU、MPU、DSP、ASIC、PLDのうちの2以上の要素の組み合わせであってもよい。なお、例えば、プロセッサ101は、図1に記載の取得部2、検出部3、算出部4、判定部5等の機能ブロックの処理を実行することが出来る。例えば、プロセッサ101は、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として実現された図1の取得部2、検出部3、算出部4、判定部5の処理を実施することができる。

【0076】

RAM102は、コンピュータ100の主記憶装置として使用される。RAM102には、プロセッサ101に実行させるOS(Operating System)のプログラムやアプリケーションプログラムの少なくとも一部が一時的に格納される。例えば、RAM102には、図1の取得部2、検出部3、算出部4、判定部5の処理を実現する機能モジュールのプログラムの一部が一時的に格納される。また、RAM102には、プロセッサ101による処理に必要な各種データ(例えば、上述の第1閾値ないし第6閾値や、算出部4が算出する相関係数など)が格納される。バス109に接続されている周辺機器としては、HDD(Hard Disk Drive)103、グラフィック処理装置104、入力インタフェース105、光学ドライブ装置106、機器接続インタフェース107およびネットワークインタフェース108がある。

【0077】

HDD103は、内蔵したディスクに対して、磁氣的にデータの書き込みおよび読み出しを行う。HDD103は、例えば、コンピュータ100の補助記憶装置として使用される。HDD103には、OSのプログラム、アプリケーションプログラム、および各種データが格納される。なお、補助記憶装置としては、フラッシュメモリなどの半導体記憶装置を使用することも出来る。なお、HDD103に、プロセッサ101による処理に必要な各種データ(例えば、上述の第1閾値ないし第6閾値や、図1の算出部4が算出する相関係数など)が格納されても良い。

【0078】

グラフィック処理装置104には、モニタ110が接続されている。グラフィック処理装置104は、プロセッサ101からの命令にしたがって、各種画像をモニタ110の画面に表示させる。モニタ110としては、CRT(Cathode Ray Tube)を用いた表示装置や液晶表示装置などがある。なお、モニタ110は、例えば、図1の判定部5の判定結果を出力する外部装置に相当する。

【0079】

入力インタフェース105には、キーボード111とマウス112とが接続されている。入力インタフェース105は、キーボード111やマウス112から送られてくる信号をプロセッサ101に送信する。なお、マウス112は、ポインティングデバイスの一例であり、他のポインティングデバイスを使用することもできる。他のポインティングデバイスとしては、タッチパネル、タブレット、タッチパッド、トラックボールなどがある。キーボード111やマウス112は、例えば、音声処理装置1を使用するユーザが、音声処理の開始や終了を指示する場合に使用されれば良い。

【0080】

光学ドライブ装置106は、レーザ光などを利用して、光ディスク113に記録されたデータの読み取りを行う。光ディスク113は、光の反射によって読み取り可能なようにデータが記録された可搬型の記録媒体である。光ディスク113には、DVD(Digital Versatile Disc)、DVD-RAM、CD-ROM(Compact Disc Read Only Memory)、CD-R(Recordable)/RW(ReWritable)などがある。可搬型の記録媒体となる光ディスク113に格納されたプログラムは光学ドライブ装置106を介して音声処理装置1にインス

10

20

30

40

50

トールされる。インストールされた所定のプログラムは、音声処理装置 1 より実行可能となる。

【0081】

機器接続インタフェース 107 は、コンピュータ 100 に周辺機器を接続するための通信インタフェースである。例えば、機器接続インタフェース 107 には、メモリ装置 114 やメモリリーダライタ 115 を接続することが出来る。メモリ装置 114 は、機器接続インタフェース 107 との通信機能を搭載した記録媒体である。メモリリーダライタ 115 は、メモリカード 116 へのデータの書き込み、またはメモリカード 116 からのデータの読み出しを行う装置である。メモリカード 116 は、カード型の記録媒体である。また、機器接続インタフェース 107 には、マイクロフォン 118 を（有線または無線通信状態）接続することができる。なお、マイクロフォン 118 は、複数接続されており、例えば、第 1 ユーザの第 1 音声や第 2 ユーザの第 2 音声が入力される。なお、第 1 音声や第 2 音声は、例えば、機器接続インタフェース 107 を介して、プロセッサ 101 に出力されることで、プロセッサ 101 は、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として実現された図 1 の取得部 2 の処理を実行することができる。

10

【0082】

ネットワークインタフェース 108 は、ネットワーク 117 に接続されている。ネットワークインタフェース 108 は、ネットワーク 117 を介して、他のコンピュータまたは通信機器との間でデータの送受信を行う。なお、ネットワークインタフェース 108 は、第 1 ユーザの第 1 音声や第 2 ユーザの第 2 音声を、ネットワーク 117 を介して他のコンピュータまたは通信機器から受け取っても良い。この場合、なお、第 1 音声や第 2 音声は、例えば、ネットワークインタフェース 108 を介して、プロセッサ 101 に出力されることで、プロセッサ 101 は、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として実現された図 1 の取得部 2 の処理を実行することができる。

20

【0083】

コンピュータ 100 は、たとえば、コンピュータ読み取り可能な記録媒体に記録されたプログラムを実行することにより、上述した音声処理機能を実現する。コンピュータ 100 に実行させる処理内容を記述したプログラムは、様々な記録媒体に記録しておくことが出来る。上記プログラムは、1 つのまたは複数の機能モジュールから構成することが出来る。例えば、図 1 に記載の取得部 2、検出部 3、算出部 4、判定部 5 等の処理を実現させた機能モジュールからプログラムを構成することが出来る。なお、コンピュータ 100 に実行させるプログラムを HDD 103 に格納しておくことができる。プロセッサ 101 は、HDD 103 内のプログラムの少なくとも一部を RAM 102 にロードし、プログラムを実行する。また、コンピュータ 100 に実行させるプログラムを、光ディスク 113、メモリ装置 114、メモリカード 116 などの可搬型記録媒体に記録しておくことも出来る。可搬型記録媒体に格納されたプログラムは、例えば、プロセッサ 101 からの制御により、HDD 103 にインストールされた後、実行可能となる。またプロセッサ 101 が、可搬型記録媒体から直接プログラムを読み出して実行することも出来る。

30

【0084】

以上に図示した各装置の各構成要素は、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。また、上記の実施例で説明した各種の処理は、予め用意されたプログラムをパーソナルコンピュータやワークステーションなどのコンピュータで実行することによって実現することができる。

40

【0085】

以上、説明した実施形態に関し、更に以下の付記を開示する。

(付記 1)

第 1 音声が含まれる第 1 入力信号と、第 2 音声が含まれる第 2 入力信号を取得する取得部と、

50

前記第 1 入力信号の第 1 信号強度と、前記第 2 入力信号の第 2 信号強度を検出する検出部と、

前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との相関係数を算出する算出部と、

前記相関係数に基づいて、前記第 1 音声と前記第 2 音声がかかっている状態か否かを判定する判定部

を備えることを特徴とする音声処理装置。

(付記 2)

前記判定部は、前記相関係数が負であり、前記相関係数が所定の第 1 閾値未満の場合に、前記第 1 音声と前記第 2 音声がかかっている状態と判定することを特徴とする付記 1 記載の音声処理装置。

10

(付記 3)

前記第 1 信号強度は前記第 1 音声のパワーまたは信号対雑音比であり、前記第 2 信号強度は前記第 2 音声のパワーまたは信号対雑音比であることを特徴とする付記 1 または付記 2 に記載の音声処理装置。

(付記 4)

前記検出部は、

前記第 1 入力信号に含まれる第 1 発話区間を前記第 1 信号強度に基づいて検出し、

前記第 2 入力信号に含まれる第 2 発話区間を前記第 2 信号強度に基づいて検出し、

前記第 1 発話区間と前記第 2 発話区間が重複する重複発話区間を検出し、

20

前記算出部は、

前記重複発話区間が所定の第 2 閾値未満の前記第 1 信号強度と前記第 2 信号強度を、前記相関係数の算出に用いないことを特徴とする付記 1 ないし付記 3 の何れか一つに記載の音声処理装置。

(付記 5)

前記検出部は、

前記第 1 入力信号に含まれる第 1 無音区間を前記第 1 信号強度に基づいて検出し、

前記第 2 入力信号に含まれる第 2 無音区間を前記第 2 信号強度に基づいて検出し、

前記第 1 無音区間と前記第 2 無音区間が重複する重複無音区間を検出し、

前記算出部は、

30

前記重複無音区間が所定の第 3 閾値未満の前記第 1 信号強度と前記第 2 信号強度を、前記相関係数の算出に用いないことを特徴とする付記 1 ないし付記 4 の何れか一つに記載の音声処理装置。

(付記 6)

前記算出部は、前記第 1 信号強度の第 1 位相または、前記第 2 信号強度の第 2 位相を所定の範囲で変化させて複数の前記相関係数を算出し、

前記判定部は、複数の前記相関係数のうち最小値となる前記相関係数に基づいて、前記第 1 音声と前記第 2 音声がかかっている状態か否かを判定することを特徴とする付記 1 ないし付記 5 の何れか一つに記載の音声処理装置。

(付記 7)

40

前記検出部は、

前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との大小関係を比較し、

前記大小関係が反転する反転回数が所定の第 4 閾値以上となる相関係数算出区間を検出し、

前記算出部は、

前記相関係数算出区間における前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との前記相関係数を算出することを特徴とする付記 1 ないし付記 6 の何れか一つに記載の音声処理装置。

(付記 8)

前記取得部は、第 3 音声が含まれる第 3 入力信号を更に取得し、

50

前記検出部は、前記第 3 入力信号の第 3 信号強度を更に検出し、

前記算出部は、前記第 1 信号強度の時系列、前記第 2 信号強度の時系列または前記第 3 信号強度の時系列のうち、2 つの信号強度の組み合わせに対する複数の前記相関係数を算出する算出部と、

前記判定部は、複数の前記相関係数のうち、前記相関係数が最小値となる 2 つの前記信号強度の組み合わせに基づいて、前記第 1 音声、前記第 2 音声または前記第 3 音声から、会話している音声の組み合わせを判定することを特徴とする付記 1 ないし付記 7 の何れか一つに記載の音声処理装置。

(付記 9)

前記算出部は、前記最小値となる前記相関係数の算出に用いた 2 つの前記信号強度を加算した加算信号強度を算出し、

10

前記加算信号強度の時系列と、前記最小値となる前記相関係数の算出に用いた 2 つの前記信号強度以外の 1 つの信号強度の時系列との参照相関係数を算出し、

前記判定部は、前記最小値となる前記相関係数が前記参照相関係数を下回る場合、または、前記第 1 閾値未満の場合に、前記参照相関係数の算出に用いた 3 つの前記信号強度に基づいて、会話している音声の組み合わせを判定することを特徴とする付記 8 記載の音声処理装置。

(付記 10)

前記算出部は、前記最小値である前記相関係数の算出に用いた 2 つの前記信号強度の前記重複無音区間が第 5 閾値以上の場合に、前記参照相関係数を算出することを特徴とする付記 9 記載の音声処理装置。

20

(付記 11)

前記検出部は、

前記第 1 入力信号に含まれる前記第 2 信号強度または、前記第 2 信号に含まれる前記第 1 信号強度を更に検出し、

前記算出部は、

前記第 1 入力信号に含まれる前記第 2 信号強度の時系列と前記第 2 入力信号に含まれる前記第 2 信号強度の時系列との第 2 相関係数、または、

前記第 2 入力信号に含まれる前記第 1 信号強度の時系列と前記第 1 入力信号に含まれる前記第 1 信号強度の時系列との第 3 相関係数を更に算出し、

30

前記判定部は、前記第 2 相関係数または前記第 3 相関係数に基づいて、前記第 1 音声と前記第 2 音声がかかっている状態か否かを判定することを特徴とする付記 1 ないし付記 10 の何れか一つに記載の音声処理装置。

(付記 12)

前記判定部は、前記相関係数が正であり、前記第 2 相関係数が所定の第 6 閾値以上の場合に、前記第 1 音声と前記第 2 音声がかかっている状態と判定することを特徴とする付記 11 記載の音声処理装置。

(付記 13)

第 1 音声が含まれる第 1 入力信号と、第 2 音声が含まれる第 2 入力信号を取得し、

前記第 1 入力信号の第 1 信号強度と、前記第 2 入力信号の第 2 信号強度を検出し、

40

前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との相関係数を算出し、

前記相関係数に基づいて、前記第 1 音声と前記第 2 音声がかかっている状態か否かを判定することを特徴とする音声処理方法。

(付記 14)

前記判定することは、前記相関係数が負であり、前記相関係数が所定の第 1 閾値未満の場合に、前記第 1 音声と前記第 2 音声がかかっている状態と判定することを特徴とする付記 13 記載の音声処理方法。

(付記 15)

前記第 1 信号強度または前記第 2 信号強度は、前記第 1 音声または前記第 2 音声のパワーまたは信号対雑音比であることを特徴とする付記 13 または付記 14 に記載の音声処理

50

方法。

(付記 16)

前記検出することは、

前記第 1 入力信号に含まれる第 1 発話区間を前記第 1 信号強度に基づいて検出し、
前記第 2 入力信号に含まれる第 2 発話区間を前記第 2 信号強度に基づいて検出し、
前記第 1 発話区間と前記第 2 発話区間が重複する重複発話区間を検出し、

前記算出部は、

前記重複発話区間が所定の第 2 閾値未満の前記第 1 信号強度と前記第 2 信号強度を、前記相関係数の算出に用いないことを特徴とする付記 13 ないし付記 15 の何れか一つに記載の音声処理方法。

10

(付記 17)

前記検出することは、

前記第 1 入力信号に含まれる第 1 無音区間を前記第 1 信号強度に基づいて検出し、
前記第 2 入力信号に含まれる第 2 無音区間を前記第 2 信号強度に基づいて検出し、
前記第 1 無音区間と前記第 2 無音区間が重複する重複無音区間を検出し、

前記算出することは、

前記重複無音区間が所定の第 3 閾値未満の前記第 1 信号強度と前記第 2 信号強度を、前記相関係数の算出に用いないことを特徴とする付記 13 ないし付記 16 の何れか一つに記載の音声処理方法。

(付記 18)

前記算出することは、前記第 1 信号強度の第 1 位相または、前記第 2 信号強度の第 2 位相を所定の範囲で変化させて複数の前記相関係数を算出し、

前記判定することは、複数の前記相関係数のうち最小値となる前記相関係数に基づいて、前記第 1 音声と前記第 2 音声がかかっている状態か否かを判定することを特徴とする付記 13 ないし付記 17 の何れか一つに記載の音声処理方法。

20

(付記 19)

前記検出することは、

前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との大小関係を比較し、

前記大小関係が反転する反転回数が所定の第 4 閾値以上となる相関係数算出区間を検出し、

30

前記算出部は、

前記相関係数算出区間における前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との前記相関係数を算出することを特徴とする付記 13 ないし付記 18 の何れか一つに記載の音声処理方法。

(付記 20)

前記取得することは、第 3 音声が含まれる第 3 入力信号を更に取得し、

前記検出することは、前記第 3 入力信号の第 3 信号強度を更に検出し、

前記算出することは、前記第 1 信号強度の時系列、前記第 2 信号強度の時系列または前記第 3 信号強度の時系列から、2 つの信号強度の組み合わせに対する複数の前記相関係数を算出する算出部と、

40

前記判定部は、複数の前記相関係数のうち、前記相関係数が最小値となる 2 つの前記信号強度の組み合わせに基づいて、前記第 1 音声、前記第 2 音声または前記第 3 音声から、かかっている音声の組み合わせを判定することを特徴とする付記 13 ないし付記 19 の何れか一つに記載の音声処理方法。

(付記 21)

前記算出することは、前記最小値となる前記相関係数の算出に用いた 2 つの前記信号強度を加算した加算信号強度を算出し、

前記加算信号強度の時系列と、前記最小値となる前記相関係数の算出に用いた 2 つの前記信号強度以外の 1 つの信号強度の時系列との参照相関係数を算出し、

前記判定することは、前記最小値となる前記相関係数が前記参照相関係数を下回る場合

50

、または、前記第 1 閾値未満の場合に、前記参照相関係数の算出に用いた 3 つの前記信号強度に基づいて、会話している音声の組み合わせを判定することを特徴とする付記 20 記載の音声処理方法。

(付記 22)

前記算出することは、前記最小値となる前記相関係数の算出に用いた 2 つの前記信号強度の前記重複無音区間が第 5 閾値以上の場合に、前記参照相関係数を算出することを特徴とする付記 21 記載の音声処理方法。

(付記 23)

前記検出することは、

前記第 1 入力信号に含まれる前記 2 信号強度または、前記第 2 入力信号に含まれる前記第 1 信号強度を更に検出し、

前記算出することは、

前記第 1 入力信号に含まれる前記第 2 信号強度の時系列と前記第 2 入力信号に含まれる前記第 2 信号強度の時系列との第 2 相関係数、または、

前記第 2 入力信号に含まれる前記第 1 信号強度の時系列と前記第 1 入力信号に含まれる前記第 1 信号強度の時系列との第 3 相関係数を更に算出し、

前記判定部は、前記第 2 相関係数または前記第 3 相関係数に基づいて、前記第 1 音声と前記第 2 音声がかかっている状態か否かを判定することを特徴とする付記 13 ないし付記 22 の何れか一つに記載の音声処理方法。

(付記 24)

前記判定することは、前記相関係数が正であり、前記第 2 相関係数が所定の第 6 閾値以上の場合に、前記第 1 音声と前記第 2 音声がかかっている状態と判定することを特徴とする付記 23 記載の音声処理方法。

(付記 25)

コンピュータに

第 1 音声が含まれる第 1 入力信号と、第 2 音声が含まれる第 2 入力信号を取得し、前記第 1 入力信号の第 1 信号強度と、前記第 2 入力信号の第 2 信号強度を検出し、前記第 1 信号強度の時系列と前記第 2 信号強度の時系列との相関係数を算出し、前記相関係数に基づいて、前記第 1 音声と前記第 2 音声がかかっている状態か否かを判定する

ことを実行させることを特徴とする音声処理プログラム。

【符号の説明】

【0086】

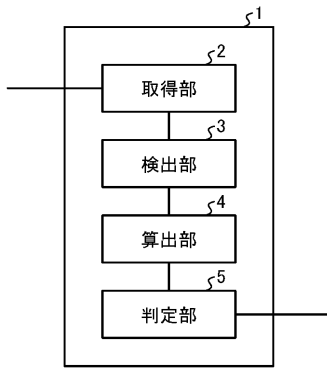
- 1 音声処理装置
- 2 取得部
- 3 検出部
- 4 算出部
- 5 判定部

10

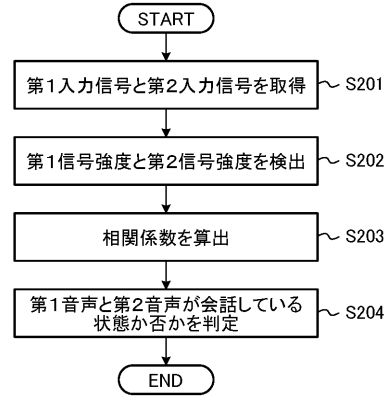
20

30

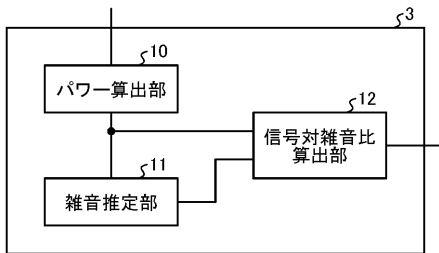
【図1】



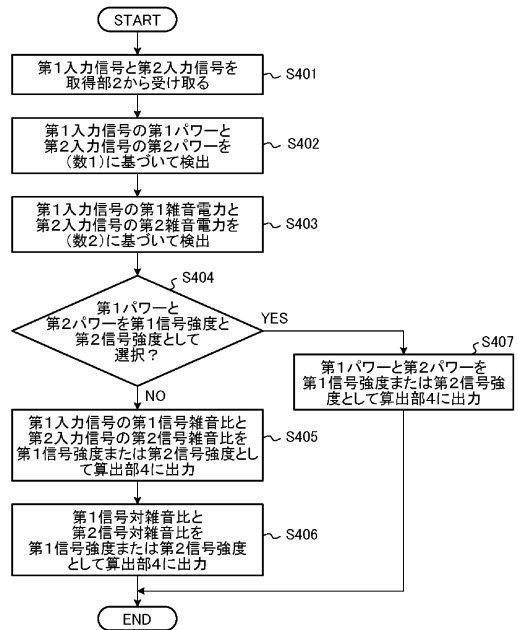
【図2】



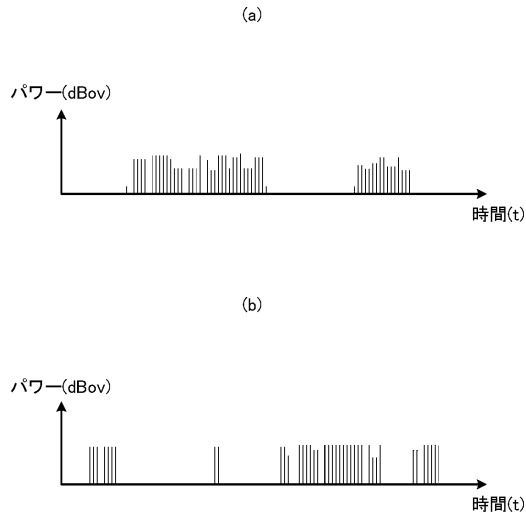
【図3】



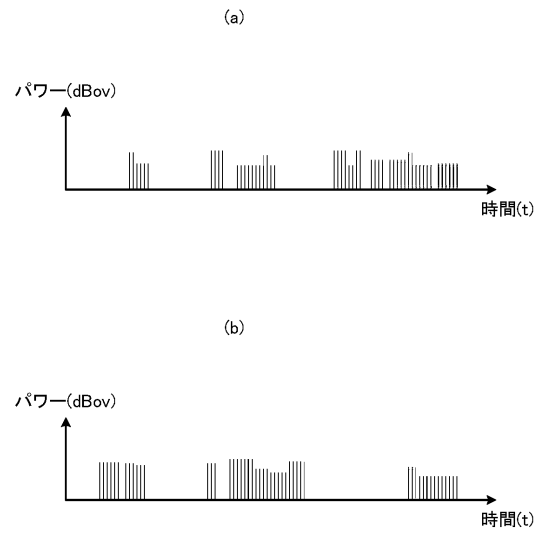
【図4】



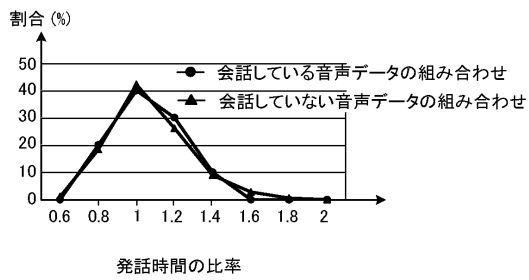
【 図 5 】



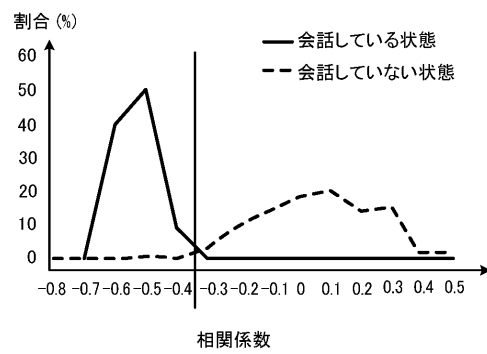
【 図 6 】



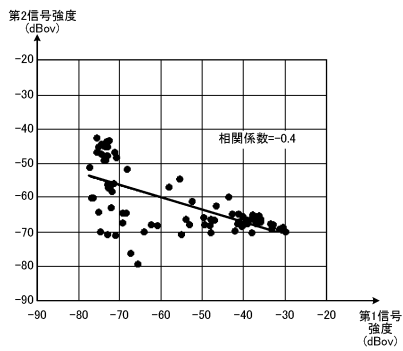
【 図 7 】



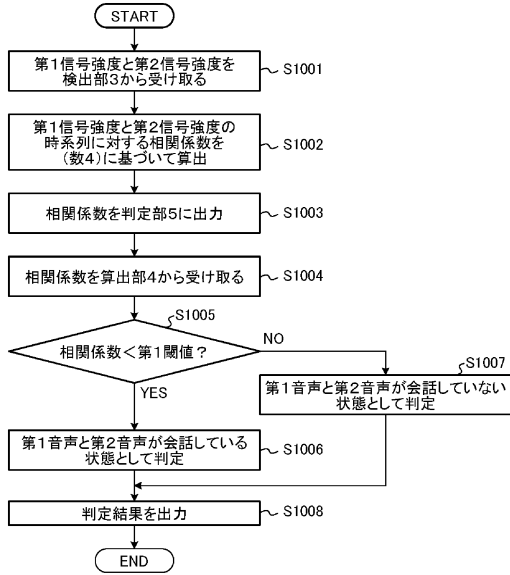
【 図 9 】



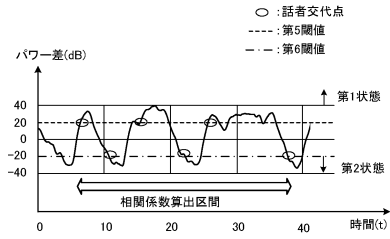
【 図 8 】



【図10】



【図12】



【図13】

相関係数算出区間に含まれる話者交替回数	会話している状態における相関係数	会話していない状態における相関係数
8 >	-0.65	0.00
8	-0.65	0.00
6	-0.60	0.04
4	-0.56	-0.60
2	-0.36	-0.45

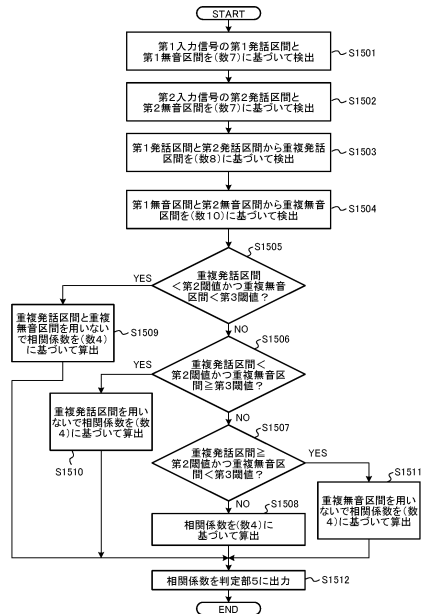
【図11】

	検出率 (%)	正解率 (%)
比較例	100	5
実施例	100	91

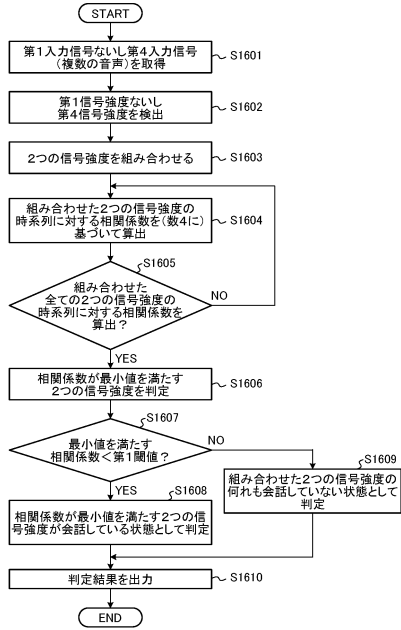
【図14】

組み合わせID	信号強度の組み合わせ	相関係数
1	第1信号強度 第2信号強度	-0.54
2	第1信号強度 第3信号強度	0.21
3	第1信号強度 第4信号強度	0.46
4	第2信号強度 第3信号強度	-0.61
5	第2信号強度 第4信号強度	-0.21
6	第3信号強度 第4信号強度	0.63

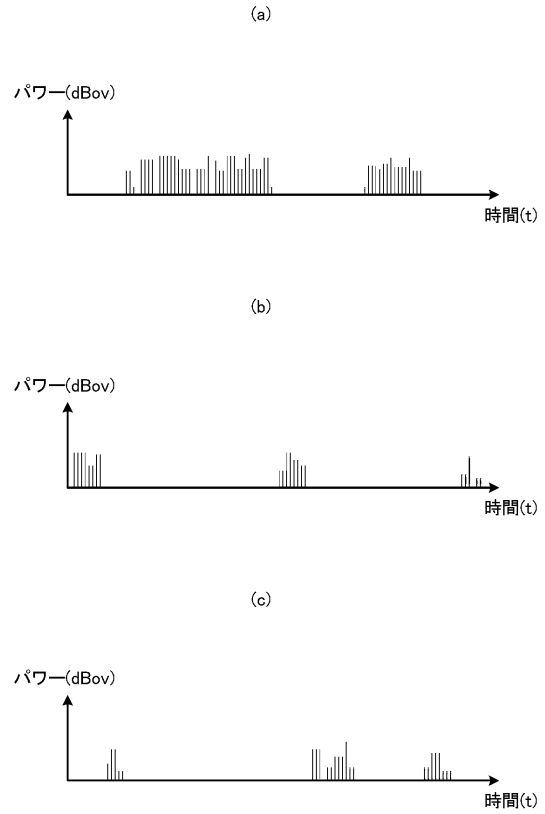
【図15】



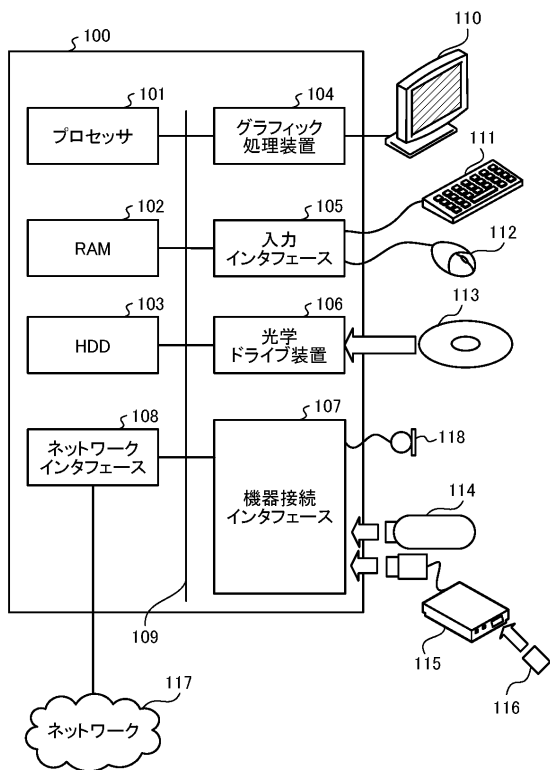
【図16】



【図17】



【図18】



フロントページの続き

審査官 山下 剛史

- (56)参考文献 特開2013-140534(JP,A)
特開2008-242318(JP,A)
特開2005-202035(JP,A)
特表2012-503400(JP,A)
特開2004-133403(JP,A)
特開2013-115622(JP,A)
Aki HARMA, et al., CONVERSATION DETECTION IN AMBIENT TELEPHONY, ICASSP 2009, IEEE, 2009年4月, p.4641-4644
川口洋平他, 音源分離を利用した話者交替の性質に基づく会話抽出, 日本音響学会2010年春季研究発表会講演論文集[CD-ROM], 2010年3月, p.869-870
- (58)調査した分野(Int.Cl., DB名)
G10L 13/00-99/00
IEEE Xplore