

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5672155号
(P5672155)

(45) 発行日 平成27年2月18日(2015.2.18)

(24) 登録日 平成27年1月9日(2015.1.9)

(51) Int.Cl.		F I			
G 1 0 L	15/04	(2013.01)	G 1 0 L	15/04	3 0 0 A
G 1 0 L	15/02	(2006.01)	G 1 0 L	15/02	2 0 0 D
G 1 0 L	25/93	(2013.01)	G 1 0 L	25/93	

請求項の数 4 (全 20 頁)

(21) 出願番号	特願2011-122808 (P2011-122808)	(73) 特許権者	000005223
(22) 出願日	平成23年5月31日(2011.5.31)		富士通株式会社
(65) 公開番号	特開2012-252060 (P2012-252060A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成24年12月20日(2012.12.20)	(74) 代理人	100089118
審査請求日	平成26年3月4日(2014.3.4)		弁理士 酒井 宏明
		(72) 発明者	張 寛
			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
		審査官	上田 雄

最終頁に続く

(54) 【発明の名称】 話者判別装置、話者判別プログラム及び話者判別方法

(57) 【特許請求の範囲】

【請求項1】

各々の話者に配置される複数のマイクから各々の音声データを取得する取得部と、前記取得部によって取得された音声データを所定の区間のフレームにフレーム化するフレーム化部と、

第1の確率モデルに基づいて、前記フレーム化部によってフレーム化されたフレームが有声音領域または無声音領域のいずれであるかを識別する第1の識別部と、

各音声データにおける同一区間のフレームで有声音領域が重複して識別された場合に、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレームの識別結果を有効化する有効化部と、

第2の確率モデルに基づいて、前記有効化部によって有効化された後のフレームの識別結果から各々の音声データにおける発話領域および沈黙領域を識別する第2の識別部とを有することを特徴とする話者判別装置。

【請求項2】

前記有効化部は、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレーム以外の識別結果を無声音領域に置き換えることを特徴とする請求項1に記載の話者判別装置。

【請求項3】

コンピュータに、

各々の話者に配置される複数のマイクから各々の音声データを取得し、

取得された音声データを所定の区間のフレームにフレーム化し、
 第1の確率モデルに基づいて、前記フレームが有声音領域または無声音領域のいずれであるかを識別し、
 各音声データにおける同一区間のフレームで有声音領域が重複して識別された場合に、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレームの識別結果を有効化し、
 第2の確率モデルに基づいて、有効化された後のフレームの識別結果から各々の音声データにおける発話領域および沈黙領域を識別する
 各処理を実行させることを特徴とする話者判別プログラム。

【請求項4】

10

コンピュータが、
 各々の話者に配置される複数のマイクから各々の音声データを取得し、
 取得された音声データを所定の区間のフレームにフレーム化し、
 第1の確率モデルに基づいて、前記フレームが有声音領域または無声音領域のいずれであるかを識別し、
 各音声データにおける同一区間のフレームで有声音領域が重複して識別された場合に、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレームの識別結果を有効化し、
 第2の確率モデルに基づいて、有効化された後のフレームの識別結果から各々の音声データにおける発話領域および沈黙領域を識別する
 各処理を実行することを特徴とする話者判別方法。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、話者判別装置、話者判別プログラム及び話者判別方法に関する。

【背景技術】

【0002】

複数の話者によってなされる会話の各場面において各話者のうち誰が発話しているのかを判別する技術が知られている。

【0003】

30

かかる話者の判別を閾値判定により実現する技術の一例として、音声認識装置が挙げられる。この音声認識装置には、各参加者に対応してマイクロホンが接続される。このような構成の下、音声認識装置は、マイクロホンによって出力される音声信号のパワーがパワー閾値を超えてから下回るまでの区間の音声信号を音声認識の対象として記憶部の所定のエリアへ記録する。その上で、音声認識装置は、記憶部に記録した音声信号を音声認識した後、発言者を特定するためのデータとしてマイクロホンの識別情報を紐付けて音声認識の結果を記憶部の議事録エリアへ記録する。

【0004】

また、話者の判別を音源定位により実現する技術の一例としては、発話イベント分離システムが挙げられる。この発話イベント分離システムでは、それぞれ異なる方向に放射状に向けた複数のマイクロホンを有するマイクロホンアレイが用いられる。発話イベント分離システムは、音源定位のアルゴリズムを用いて、マイクロホンアレイによって収録された多チャンネルの音声データを解析して時刻毎に音の到来方向を推定する。また、発話イベント分離システムは、音源となる話者の存在範囲を推定する。その上で、発話イベント分離システムは、音源定位の結果と、話者の存在範囲の推定結果から、時刻毎にどの話者が発話しているかを同定する。

40

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2008-309856号公報

50

【特許文献2】特開2007-233239号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかしながら、上記の従来技術では、以下に説明するように、話者の判別を簡易かつ正確に行うことができないという問題がある。

【0007】

例えば、上記の音声認識装置は、音声信号のパワーがパワー閾値を超過するか否かによって話者が発話しているか否かを判定するものである。このため、上記の音声認識装置では、話者を判別する精度はパワー閾値に依存するが、人間が発話する音声には個人差がある

10

【0008】

また、上記の発話イベント分離システムでは、音源定位により音の到来方向を推定するのに複雑なアルゴリズムを使用する必要がある。さらに、上記の発話イベント分離システムでは、話者の存在範囲を推定するために、会議に参加する人数等を予め学習させておく必要もある。よって、上記の発話イベント分離システムでは、話者の判別を簡易に行うことはできない。

【0009】

開示の技術は、上記に鑑みてなされたものであって、話者の判別を簡易かつ正確に行うことができる話者判別装置、話者判別プログラム及び話者判別方法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】

【0010】

本願の開示する話者判別装置は、各々の話者に配置される複数のマイクから各々の音声データを取得する取得部を有する。さらに、前記話者判別装置は、前記取得部によって取得された音声データを所定の区間のフレームにフレーム化するフレーム化部を有する。さらに、前記話者判別装置は、第1の確率モデルに基づいて、前記フレーム化部によってフレーム化されたフレームが有声音領域または無声音領域のいずれであるかを識別する第1の識別部を有する。さらに、前記話者判別装置は、各音声データにおける同一区間のフレームで有声音領域が重複して識別された場合に、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレームの識別結果を有効化する有効化部を有する。さらに、前記話者判別装置は、第2の確率モデルに基づいて、前記有効化部によって有効化された後のフレームの識別結果から各々の音声データにおける発話領域および沈黙領域を識別する第2の識別部を有する。

30

【発明の効果】

【0011】

本願の開示する話者判別装置の一つの態様によれば、話者の判別を簡易かつ正確に行うことができるという効果を奏する。

【図面の簡単な説明】

40

【0012】

【図1】図1は、実施例1に係る会話分析装置の機能的構成を示すブロック図である。

【図2】図2は、有声音および無声音の一例を示す図である。

【図3】図3は、発話領域および沈黙領域の一例を示す図である。

【図4】図4は、話者判別方法を説明するための図である。

【図5】図5は、隠れマルコフモデルにおける状態遷移図の一例を示す図である。

【図6】図6は、有声音領域および無声音領域の識別結果の一例を示す図である。

【図7】図7は、図6に示した識別結果の置換結果の一例を示す図である。

【図8】図8は、隠れマルコフモデルにおける状態遷移図の一例を示す図である。

【図9】図9は、実施例1に係る会話分析処理の手順を示すフローチャートである。

50

【図10】図10は、実施例1に係る会話分析処理の手順を示すフローチャートである。

【図11】図11は、実施例1に係る有効化処理の手順を示すフローチャートである。

【図12】図12は、実施例1及び実施例2に係る話者判別プログラムを実行するコンピュータの一例について説明するための図である。

【発明を実施するための形態】

【0013】

以下に、本願の開示する話者判別装置、話者判別プログラム及び話者判別方法の実施例を図面に基づいて詳細に説明する。なお、この実施例は開示の技術を限定するものではない。そして、各実施例は、処理内容を矛盾させない範囲で適宜組み合わせることが可能である。

10

【実施例1】

【0014】

まず、本実施例に係る話者判別装置を含む会話分析装置の機能的構成について説明する。図1は、実施例1に係る会話分析装置の機能的構成を示すブロック図である。図1に示す会話分析装置10は、話者A、話者B及び話者Cにそれぞれ対応して設けられた接話マイク30A～30Cを介して集音した複数の音声データから、話者A～話者Cの会話に関する特性を抽出して会話スタイルを分析するものである。

【0015】

この会話分析装置10には、接話マイク30A～30Cの3つのマイクが接続される。これら接話マイク30A～30Cは、話者によって装着される接話型マイクロホン(close talking microphone)である。かかる接話マイクの一態様としては、ラベルマイクやヘッドセットマイクなどが挙げられる。以下では、接話マイク30A～30Cのことを区別なく総称する場合には「接話マイク30」と記載する場合がある。

20

【0016】

なお、図1の例では、接話型マイクロホンを用いる場合を例示したが、必ずしも接話型マイクロホンを用いる必要はなく、各々の話者に他の話者よりも接近して配置するのであれば任意のマイクを採用できる。また、図1の例では、3つのマイクを用いて話者A～話者Cの3人の会話を集音する場合を例示するが、2つのマイクを用いて2人の会話を集音することとしてもよいし、また、4つ以上のマイクを用いて4人以上の会話を集音することとしてもかまわない。

30

【0017】

登録部31は、接話マイク30によって集音された音声信号を会話分析装置10の記憶部11へ登録する処理部である。一態様としては、登録部31は、接話マイク30から音声入力されたアナログ信号にA/D(Analog/Digital)変換を実行することによりデジタル信号に変換した上で音声記憶部11へ登録する。なお、以下では、接話マイク30Aから音声入力されたアナログ信号がA/D変換されたデジタル信号のことを「第1の音声データ」と記載する場合がある。また、接話マイク30Bから音声入力されたアナログ信号がA/D変換されたデジタル信号のことを「第2の音声データ」と記載する場合がある。さらに、接話マイク30Cから音声入力されたアナログ信号がA/D変換されたデジタル信号のことを「第3の音声データ」と記載する場合がある。

40

【0018】

図1に示すように、会話分析装置10は、音声記憶部11と、抽出部13と、分析部14とを有する。なお、会話分析装置10は、図1に示した機能部以外にも既知のコンピュータが有する各種の機能部、例えば各種の入力デバイスや音声出力デバイスなどを始め、他の装置との通信を制御する通信インターフェースなどの機能部を有するものとする。

【0019】

音声記憶部11は、音声データを記憶する記憶部である。この音声記憶部11は、第1の音声データ12Aと、第2の音声データ12Bと、第3の音声データ12Cとを記憶する。

【0020】

50

これら第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cは、話者A～話者Cが装着する接話マイク30によって集音された音声信号がA/D変換されたデジタルデータである。このうち、第1の音声データ12Aには、話者Aの音声だけでなく、話者Bおよび話者Cの音声も含み得るが、話者Aから接話マイク30Aまでの距離が話者Bや話者Cに比べて接近している。よって、第1の音声データ12Aに含まれる音声は、話者Aと話者Bや話者Cとの間で同時に発話が行なわれていた場合でも、話者Aによって発話された音声のエネルギーが最も高くなる。同様に、第2の音声データ12Bに含まれる音声は、話者Bによって発話された音声のエネルギーが最も高くなり、第3の音声データ12Cに含まれる音声は、話者Cによって発話された音声のエネルギーが最も高くなる。

10

【0021】

なお、上記の音声記憶部11などの記憶部には、半導体メモリ素子や記憶装置を採用できる。例えば、半導体メモリ素子としては、V R A M (Video Random Access Memory)、R A M (Random Access Memory)、R O M (Read Only Memory)やフラッシュメモリ (flash memory)などが挙げられる。また、記憶装置としては、ハードディスク、光ディスクなどの記憶装置が挙げられる。

【0022】

ここで、話者によって発話される有声音および無声音について説明する。図2は、有声音および無声音の一例を示す図である。図2の例では、サンプリング周波数が16kHzである接話マイクを用いて取得した音声データが示されている。図2の例では、横軸は時間を示し、縦軸は周波数を示し、図中の濃淡はスペクトルエントロピーの大小を示す。

20

【0023】

図2に示すように、有声音V (Voiced)は、スペクトルエントロピーの変化が大きく、無声音U (Unvoiced)よりも低い周波数の音である。有声音の一例としては、母音「a」、「i」、「u」、「e」、「o」などが挙げられる。また、無声音Uは、有声音Vよりも高い周波数の音である。無声音の一例としては、母音以外の音、例えば「s」、「p」、「h」などが挙げられる。これら有声音および無声音の特徴は、話者によって発話される言語に依存せず、日本語、英語や中国語などの任意の言語において共通する。

【0024】

次に、有声音および無声音と発話領域および沈黙領域との関係について説明する。図3は、発話領域および沈黙領域の一例を示す図である。発話領域は、話者によって発話が行なわれている領域を指し、無声音領域および有声音領域を含む。なお、図3の例では、話者によって「WaTaShiWa Chou DeSu」と発話された場合を示す。

30

【0025】

図3に示す発話の例では、「WaTaShiWa」の発話領域40と、「Chou」の発話領域41と、「DeSu」の発話領域42との間に、沈黙領域43および沈黙領域44が存在することを示す。このうち、発話領域40には、無声音「W」、有声音「a」、無声音「T」、有声音「a」、無声音「Sh」、有声音「i」、無声音「W」、有声音「a」が含まれる。また、発話領域41には、無声音「Ch」、有声音「ou」が含まれる。さらに、発話領域42には、無声音「D」、有声音「e」、無声音「S」、有声音「u」が含まれる。

40

【0026】

図1の説明に戻り、会話分析装置10は、複数の話者によってなされる会話の各場面において各話者のうち誰が発話しているのかを判別する話者判別装置50を有する。

【0027】

ここで、本実施例に係る話者判別装置50は、接話マイク30A～30Cから第1の音声データ、第2の音声データ及び第3の音声データを取得する。さらに、本実施例に係る話者判別装置50は、第1の音声データ、第2の音声データ及び第3の音声データを所定の区間のフレームにフレーム化する。さらに、本実施例に係る話者判別装置50は、第1の確率モデルに基づいて、フレームが有声音領域または無声音領域のいずれであることを識

50

別する。さらに、本実施例に係る話者判別装置50は、各音声データにおける同一区間のフレームで有声音領域が重複して識別された場合に、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレームの識別結果を有効化する。その上で、本実施例に係る話者判別装置50は、第2の確率モデルに基づいて、有効化された後のフレームの識別結果から各々の音声データにおける発話領域および沈黙領域を識別する。

【0028】

図4を用いて、上記の話者判別方法について説明する。図4は、話者判別方法を説明するための図である。図4の上段には、各フレームの有声音領域または無声音領域の識別結果が図示されている。図4の中段には、最大のエネルギーを持つフレームの有声音領域の識別結果が有効化された後の各フレームの識別結果が図示されている。図4の下段には、各々の音声データにおける発話領域および沈黙領域の識別結果が図示されている。

10

【0029】

図4の上段に示すように、話者判別装置50は、第1の確率モデルに基づいて、第1の音声データ、第2の音声データ及び第3の音声データからフレーム化した各フレームが有声音領域または無声音領域のいずれであるかを識別する。ここで、図4の例では、記号「 \square 」、記号「 \square 」、記号「 \square 」がそれぞれ音声データのフレームを表し、記号「 \square 」及び記号「 \square 」が有声音領域であることを示し、記号「 \square 」が無声音領域であることを示す。図4に示す記号「 \square 」のフレームは、図4に示す記号「 \square 」のフレームよりも高いエネルギーを有することを示す。これら第1の音声データ、第2の音声データおよび第3の音声データの識別結果からは、話者A～話者Cのうち話者Bと話者Cが会話しており、話者Bが話者Cよりも大声で発話していることが推定できる。なお、以下では、第1の音声データから得られた各フレームのことを観測順に第1フレーム(1)・・・第1フレーム(n)と記載する場合がある。また、第2の音声データから得られた各フレームのことを観測順に第2フレーム(1)・・・第2フレーム(m)と記載する場合がある。さらに、第3の音声データから得られた各フレームのことを観測順に第3フレーム(1)・・・第3フレーム(m)と記載する場合がある。

20

【0030】

また、図4の中段に示すように、話者判別装置50は、各音声データの同一の区間のフレームで有声音領域が重複する場合に、最大のエネルギーを持つフレームの識別結果を有効化する。この例では、第2の音声データ及び第3の音声データを構成するフレームのうち、下記のように、同一区間のフレームで互いに識別結果が有声音領域と識別されている。すなわち、第2フレーム(1)と第3フレーム(1)、第2フレーム(6)と第3フレーム(6)、第2フレーム(10)と第3フレーム(10)において互いの識別結果が有声音領域と識別されている。さらに、第2フレーム(13)と第3フレーム(13)、第2フレーム(18)と第3フレーム(18)において互いの識別結果が有声音領域と識別されている。この場合には、いずれのフレームについても第2の音声データのエネルギーの方が高いので、第3フレーム(1)、第3フレーム(6)、第3フレーム(10)、第3フレーム(13)及び第3フレーム(18)の識別結果が有声音から無声音に置き換えられる。

30

【0031】

さらに、図4の下段に示すように、話者判別装置50は、第2の確率モデルに基づいて、有効化後のフレームの識別結果から各々の音声データにおける発話領域および沈黙領域を識別する。この例では、第2の音声データのフレームのうち下線が引かれた領域が話者Bの発話領域として識別されている。さらに、第3の音声データのフレームのうち下線が引かれた領域が話者Cの発話領域として識別されている。この場合には、話者Bの発話領域と話者Cの発話領域が重複するフレーム、すなわち第2フレーム(7)～第2フレーム(13)の区間が同時発話として判別される。

40

【0032】

このように、本実施例に係る話者判別装置50は、各音声データにおける同一区間のフレームで有声音領域が重複して識別された場合に、最大のエネルギーを持つフレームの識

50

別結果だけを有効化して各々の音声データの発話領域および沈黙領域を識別する。このため、本実施例に係る話者判別装置50は、各音声データを構成する同一区間のフレーム間で閾値を用いて判定せずとも、話者を判別することができる。さらに、本実施例に係る話者判別装置50では、話者の判別に複雑なアルゴリズムを用いる必要はなく、事前に学習を行う必要もない。したがって、本実施例に係る話者判別装置50によれば、話者の判別を簡易かつ正確に行うことができる。

【0033】

また、本実施例に係る話者判別装置50は、各音声データにおける同一区間のフレームで有声音領域が単独で識別された場合には、エネルギーの大小に関係なく、有声音領域と識別された識別結果を維持する。一般に、発話は、有声音と無声音が混在して構成されるので、複数の話者によって同時に発話された場合でも、同時発話で有声音領域が完全に重複する可能性は低く、有声音領域が単独で識別される機会が残る可能性は高い。例えば、図4の下段の例で言えば、話者Cの発話の音量が話者Bの発話の音量よりも低くても、第3フレーム(7)、第3フレーム(9)及び第3フレーム(12)の識別結果は有声音のまま維持される。それゆえ、本実施例に係る話者判別装置50では、話者が発話する音量に開きがある場合でも、同時発話を判別することもできる。

【0034】

さらに、話者判別装置50を詳細に説明する。図1に示すように、話者判別装置50は、取得部51と、フレーム化部52と、第1の識別部53と、有効化部54と、第2の識別部55とを有する。

【0035】

取得部51は、第1の音声データ、第2の音声データおよび第3の音声データを取得する処理部である。一態様としては、取得部51は、音声記憶部11に記憶された第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cを読み出す。他の一態様としては、取得部51は、登録部31によってA/D変換された第1の音声データ、第2の音声データおよび第3の音声データをストリームデータとして取得することもできる。更なる一態様としては、取得部51は、ネットワークを介して図示しない外部装置から第1の音声データ、第2の音声データおよび第3の音声データを取得することもできる。

【0036】

フレーム化部52は、取得部51によって取得された第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cを所定の区間のフレームにフレーム化する処理部である。一態様としては、フレーム化部52は、第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cそれぞれの長さを比較する。そして、フレーム化部52は、第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cの長さの差が許容誤差範囲内でない場合には、図示しない表示部等にエラーメッセージを出力し、以降の処理を行わない。一方、フレーム化部52は、第1の音声データ12A、第2の音声データ12B及び第3の音声データの長さが同一であるか、あるいは許容誤差範囲内である場合には、下記のような処理を実行する。すなわち、フレーム化部52は、第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cをフレーム化する。一例を挙げれば、フレーム化部52は、下記の式(1)、式(2)を用いて、各々の音声データを、長さを256msとするフレーム化を行う。このとき、フレーム化部52は、前後のフレームの重複部分の長さが128msとなるようにする。なお、上記のフレームの長さ、前後のフレームの重複部分の長さは、あくまでも一例であり、任意の値を採用できる。

$$S = \text{floor}(Y/X) \dots \dots \dots \text{式(1)}$$

$$m = \text{floor}((S - 256) / 128) + 1 \dots \dots \dots \text{式(2)}$$

なお、「 $\text{floor}(x)$ 」は、 x 以下の最大の整数を算出するための関数であり、 Y は、第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cそれぞれのデータ量(byte)であり、 X は、1(byte)のデータに対応する長さ(ms)である。

10

20

30

40

50

【 0 0 3 7 】

第1の識別部53は、第1の確率モデルに基づいて、フレーム化部52によってフレーム化されたフレームが有声音領域または無声音領域のいずれであるかを識別する処理部である。一態様としては、第1の識別部53は、第1フレーム(1)～第1フレーム(m)、第2フレーム(1)～第2フレーム(m)、第3フレーム(1)～第3フレーム(m)の各々の音声データごとに、下記の処理を実行する。すなわち、第1の識別部53は、自己相関係数のピークの数、自己相関係数のピークの最大値及びスペクトルエントロピーの3つの特徴量を抽出する。さらに、第1の識別部53は、先に抽出した3つの特徴量それぞれの平均値および標準偏差を各々の音声データごとに算出する。その上で、第1の識別部53は、確率モデルである隠れマルコフモデル(Hidden Markov Model; HMM)を用いて、有声音領域および無声音領域を各々の音声データごとに識別する。

10

【 0 0 3 8 】

ここで、有声音領域および無声音領域の識別方法について説明する。図5は、隠れマルコフモデルにおける状態遷移図の一例を示す図である。図5に示すように、第1の識別部53は、上記の3つの特徴量、並びに、各特徴量の平均値および標準偏差を観測結果(observation)とし、EM法(Expectation-Maximization algorithm)を用いて、状態遷移確率(transition possibility) P_t を算出する。

【 0 0 3 9 】

かかる状態遷移確率 P_t は、例えば、有声音の状態のままの確率、有声音の状態から無声音の状態に遷移する確率、無声音の状態のままの確率、無声音の状態から有声音の状態に遷移する確率を指す。図5に示す例で言えば、発話は、有声音および無声音の両方とも同一の確率で開始すると仮定して、発話の開始における有声音および無声音の状態の確率がいずれも「0.5」と設定されている。さらに、初期の状態遷移確率 P_t として、有声音の状態のままの確率が「0.95」に設定されるとともに、有声音の状態から無声音の状態に遷移する確率が「0.05」に設定されている。さらに、初期の状態遷移確率 P_t として、無声音の状態のままの確率が「0.95」に設定されるとともに、無声音の状態から有声音の状態に遷移する確率が「0.05」に設定されている。このような設定の下、第1の識別部53は、状態遷移確率 P_t を算出することを所定回数繰り返す。これによって、精度の高い状態遷移確率 P_t を算出することができる。

20

【 0 0 4 0 】

さらに、第1の識別部53は、上記の3つの特徴量、並びに、各特徴量の平均値および標準偏差を観測結果とし、ビタビアルゴリズム(Viterbi algorithm)により、観測確率(observation possibility) P_o を各々の音声データごとに算出する。ここで、観測確率 P_o は、例えば、有声音の状態から観測(observed)を出力する確率、有声音の状態から非観測(not observed)を出力する確率、無声音の状態から観測を出力する確率および無声音の状態から非観測を出力する確率である。なお、観測確率は、出力確率(emission possibility)とも称される。

30

【 0 0 4 1 】

これら状態遷移確率 P_t および観測確率 P_o を算出した後に、第1の識別部53は、上記の3つの特徴量に基づいて、ビタビアルゴリズムを用いて、次のような処理を実行する。すなわち、第1の識別部53は、発話が行われている各フレームにおいて発話されている音が有声音であるか、あるいは無声音であるかを識別する。その上で、第1の識別部53は、有声音と識別された領域を有声音領域とし、無声音と識別された領域を無声音領域とする。

40

【 0 0 4 2 】

このように、第1の識別部53は、自己相関係数のピークの数、自己相関係数のピークの最大値及びスペクトルエントロピーなどの特徴量を用いて、有声音領域および無声音領域を識別する。したがって、第1の識別部53では、周囲のノイズの影響によって有声音領域および無声音領域を識別する精度が低下することを抑制できる。また、第1の識別部53は、周囲のノイズに強い特徴量を用いるため、第1の音声データ12A、第2の音声

50

データ12B及び第3の音声データ12Cをフレーム化する場合に、フレームの個数をより少なくすることができる。それゆえ、第1の識別部53では、より簡易な処理で有声音領域および無声音領域を識別できる。

【0043】

有効化部54は、各音声データにおける同一区間のフレームで有声音領域が重複する場合に、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレームの識別結果を有効化する。

【0044】

一態様としては、有効化部54は、各音声データにおける同一区間のフレームで第1の識別部53による識別結果を比較する。このとき、有効化部54は、同一区間のフレームで有声音領域が重複する場合に、当該有声音領域と識別されたフレームのエネルギーを演算する。そして、有効化部54は、当該同一区間で有声音領域と識別されたフレームのうち最大エネルギーを持つフレームを特定する。その上で、有効化部54は、最大エネルギーを持つフレーム以外の識別結果を有声音領域から無声音領域に置き換える。その後、有効化部54は、各音声データ間で同一区間のフレームを全て処理するまで、識別結果の比較、フレームの特定、識別結果の置き換えを繰り返し実行する。なお、上記のエネルギーは、各々の音声データのフレームに高速フーリエ変換、いわゆるFFT(Fast Fourier Transform)を実行して周波数解析を行った上で周波数成分ごとの振幅値を平均化することにより算出される。

【0045】

ここで、有効化部54による識別結果の置換要領について説明する。図6は、有声音領域および無声音領域の識別結果の一例を示す図である。図7は、図6に示した識別結果の置換結果の一例を示す図である。図6に示すように、「12時00分00.000秒」から「12時00分00.010秒」までの区間では、第1フレーム、第2フレーム及び第3フレームの全ての識別結果が有声音領域と識別されている。この場合には、有効化部54は、図7に示すように、第1フレーム、第2フレーム及び第3フレームのうちエネルギーが最大である第1フレームを除き、第2フレーム及び第3フレームの識別結果を有声音領域「V」から無声音領域「U」へ置き換える。また、図6に示す「12時00分00.010秒」から「12時00分00.020秒」までの区間では、第1フレーム及び第2フレームの識別結果が有声音領域と識別されている。この場合には、有効化部54は、図7に示すように、第1フレーム及び第2フレームのうちエネルギーが最大である第1フレームの識別結果を維持する一方で、最大でない第2フレームの識別結果を有声音領域「V」から無声音領域「U」へ置き換える。さらに、図6に示すように、「12時00分00.020秒」から「12時00分00.030秒」までの区間では、第2フレームの識別結果だけが有声音領域と識別されている。この場合には、有効化部54は、同一区間のフレームで有声音領域が重複しないので、図7に示すように、第2フレームの識別結果を維持する。

【0046】

第2の識別部55は、第2の確率モデルに基づいて、有効化部54による有効化がなされた後のフレームの識別結果から各々の音声データにおける発話領域および沈黙領域を識別する処理部である。

【0047】

ここで、発話領域および沈黙領域の識別方法について説明する。図8は、隠れマルコフモデルにおける状態遷移図の一例を示す図である。図8に示す状態遷移確率 P_{ij} および観測確率 P_o は、予め定められた値である。かかる状態遷移確率 P_{ij} は、例えば、沈黙の状態である沈黙状態のままの確率、沈黙状態から発話の状態である発話状態に遷移する確率、発話状態のままの確率および発話状態から沈黙状態に遷移する確率を示す。図8に示す例で言えば、発話は、有声音および無声音の両方とも同一の確率で開始すると仮定して、発話の開始における沈黙状態および発話状態の確率がいずれも「0.5」に設定されている。また、状態遷移確率 P_{ij} として、沈黙状態のままの確率が「0.999

10

20

30

40

50

」に設定されるとともに、沈黙状態から発話状態に遷移する確率が「0.001」に設定されている。さらに、状態遷移確率 P_1 として、発話状態のままの確率が「0.999」設定されるとともに、発話状態から沈黙状態に遷移する確率が「0.001」に設定されている。

【0048】

また、観測確率 P_0 は、例えば、沈黙状態において無声音が検出される確率、沈黙状態において有声音が検出される確率、発話状態において無声音が検出される確率、および発話状態において有声音が検出される確率を指す。図8の例で言えば、観測確率 P_0 として、沈黙状態において無声音が検出される確率が「0.99」に設定されるとともに、沈黙状態において有声音が検出される確率が「0.01」に設定されている。また、観測確率 P_0 として、発話状態において無声音が検出される確率が「0.5」に設定されるとともに、発話状態において有声音が検出される確率が「0.5」に設定されている。

10

【0049】

なお、図8の例では、発話状態において無声音が検出される確率および発話状態において有声音が検出される確率をともに「0.5」に設定する場合を例示したが、同時発話の場合には他の話者よりも音量が小さい発話を行う話者の無声音が増加することも想定される。よって、発話状態において無声音が検出される確率を「0.5」よりも大きく設定することにより、他の話者よりも音量が小さい発話を行う話者の無声音の増加を抑制することもできる。

【0050】

20

このような設定の下、第2の識別部55は、ピタビアルゴリズムを用いて、有効化部54による有効化がなされた後の有声音および無声音から、各々の音声データにおける沈黙領域および発話領域であるかを識別する。これによって、第1の音声データにおける話者Aの発話領域および沈黙領域、第2の音声データにおける話者Bの発話領域および沈黙領域、さらには、第3の音声データにおける話者Cの発話領域および沈黙領域が識別される。

【0051】

会話分析装置10の説明に戻り、抽出部13は、各々の音声データから会話特性を抽出する処理部である。一態様としては、抽出部13は、第2の識別部55によって識別された第1の音声データにおける話者Aの発話領域をもとに有声音領域の数、有声音領域の長さの平均値および有声音領域の長さの標準偏差を算出する。また、抽出部13は、第2の識別部55によって識別された第1の音声データにおける話者Aの発話領域をもとに発話領域の数、発話領域の長さの平均値および発話領域の長さの標準偏差を算出する。さらに、抽出部13は、第2の識別部55によって識別された第1の音声データにおける話者Aの沈黙領域をもとに、沈黙領域の数、沈黙領域の長さの平均値および沈黙領域の長さの標準偏差を算出する。

30

【0052】

また、抽出部13は、会話全体の時間の長さに対する話者Aの発話時間の長さの割合を算出する。このとき、抽出部13は、話者Aの発話領域の長さの合計を、話者Aの発話時間の長さとして、上記の割合を算出する。また、抽出部13は、話者Bの発話時間に対する話者Aの発話時間の割合を算出する。さらに、抽出部13は、話者Cの発話時間に対する話者Aの発話時間の割合も算出する。また、抽出部13は、話者Aの発話領域をもとに、音量の標準偏差およびスペクトルエントロピーの標準偏差を算出する。さらに、抽出部13は、話者Aの発話領域をもとに算出した音量の標準偏差と、スペクトルエントロピーの標準偏差との和を、変化の度合いとして算出する。なお、ここでは、話者Aの会話特性を抽出する場合を例示したが、話者Bおよび話者Cについても、上記の話者Aと同様にして、会話特性を抽出する。

40

【0053】

このようにして算出された有声音領域の数、有声音領域の長さの平均値および有声音領域の長さの標準偏差の各会話特性は、有声音の長さがどの位長いのかを示す指標となる。

50

また、発話領域の数、発話領域の長さの平均値、および発話領域の長さの標準偏差の各会話特性は、対応する人物が、常に会話において長く続けて話すのか、あるいは少ししか話さないのかを示す指標となる。また、沈黙領域の数、沈黙領域の長さの平均値および沈黙領域の長さの標準偏差の各会話特性は、話者の話し方が、長く続けて話すのか、あるいは中断（沈黙）を多くはさみながら話すのかを示す指標となる。また、会話全体の時間の長さに対するある人物の発話時間の長さの割合および他の人物の発話時間に対するある人物の発話時間の割合 R_t の各会話特性は、会話の参加状態を示す指標となる。また、音量の標準偏差、スペクトルエントロピーの標準偏差および変化の度合いの各会話特性は、感情の変化が激しい情熱的な話者であるのか、あるいは感情の変化が小さい静かな話者であるのかを示す指標となる。

10

【0054】

分析部14は、抽出部13によって抽出された会話特性に基づいて、会話スタイルを分析する処理部である。一態様としては、分析部14は、他の人物の発話時間に対するある人物の発話時間の割合 R_t が、所定値、例えば1.5以上である場合には、この「ある人物」は、会話においてよく話す人物であると分析する。また、分析部14は、割合 R_t が所定値、例えば0.66以下である場合には、この「ある人物」は、会話においてあまり話さない、いわゆる聞き役の人物であると分析する。なお、分析部14は、割合 R_t が、所定値、例えば0.66より大きく、1.5未満である場合には、会話に対する参加状況において両者は対等であると分析する。

【0055】

20

他の一態様としては、分析部14は、ある人物の発話領域の数に対する有声音領域の数の割合および発話領域の長さの平均値が、他の人物の発話領域の数に対する有声音領域の数の割合および発話領域の長さの平均値よりも大きい場合には、次のように分析する。すなわち、分析部14は、「ある人物」は会話において長く続けて話しがちな人物であると分析する。また、分析部14は、ある人物の沈黙領域の長さの平均値が他の人物の沈黙領域の長さの平均値よりも大きく、かつある人物の沈黙領域の長さの標準偏差が所定値、例えば、6.0以上である場合には、次のように分析する。すなわち、分析部14は、「ある人物」は、相手の話を聞いて、相手の内容に合わせて自分の発話を中断するため、発話の長さが一定しない人物であると分析する。

【0056】

30

更なる一態様としては、分析部14は、ある人物の音量の標準偏差、スペクトルエントロピーの標準偏差または変化の度合いが、それぞれに対応する基準値以上である場合には、「ある人物」は感情の変化が激しい情熱的な話者であると分析する。また、分析部14は、ある人物の音量の標準偏差、スペクトルエントロピーの標準偏差または変化の度合いが、それぞれに対応する基準値未満である場合には、「ある人物」は感情の変化が小さい静かな話者であると分析する。

【0057】

他の一態様としては、分析部14は、ある人物と他の人物との関係を分析することもできる。例えば、分析部14は、他の人物の発話時間に対するある人物の発話時間の割合 R_t が所定値、例えば1.0以上である場合には、「ある人物」は「他の人物」に対してよく話しかけているため、ある人物と他の人物との関係が友達や家族であると分析できる。一方、割合 R_t が所定値、例えば1.0未満である場合には、この「ある人物」は「他の人物」の話を聞こうとしているため、ある人物と他の人物との関係が会社の同僚やビジネスパートナーであると分析できる。

40

【0058】

更なる一態様としては、分析部14は、ある人物と他の人物との会話においてある人物の発話領域の長さの平均値が所定値、例えば、1.85(s)以上である場合には、ある人物と他の人物との関係が友達や家族であると分析できる。これは、「ある人物」が「他の人物」に対してよく話しかけているためである。一方、分析部14は、ある人物と他の人物との会話においてある人物の発話領域の長さの平均値が所定値、例えば、1.85(

50

s) 未満である場合には、ある人物と他の人物との関係が会社の同僚やビジネスパートナーであると分析できる。

【0059】

他の一態様としては、分析部14は、ある人物と他の人物との会話においてある人物の沈黙領域の長さの平均値が所定値、例えば、3.00(s)以下である場合には、同様の理由で、ある人物と他の人物との関係が友達や家族であると分析できる。一方、分析部14は、ある人物の沈黙領域の長さの平均値が所定値、例えば、3.00(s)より大きい場合には、ある人物と他の人物との関係が会社の同僚やビジネスパートナーであると分析できる。

【0060】

更なる一態様としては、分析部14は、ある人物と他の人物との会話においてある人物の変化の度合いが所定値、例えば、0.33以上である場合には、同様の理由で、ある人物と他の人物との関係が友達や家族であると分析できる。一方、分析部14は、ある人物の変化の度合いが所定値、例えば、0.33未満である場合には、ある人物と他の人物との関係が会社の同僚やビジネスパートナーであると分析できる。

【0061】

これらの分析を行った後に、分析部14は、分析結果を所定の出力先の装置、例えば会話分析装置10が有する表示部や話者A～話者Cが利用する情報処理装置などに出力することができる。

【0062】

なお、話者判別装置50、抽出部13及び分析部14には、各種の集積回路や電子回路を採用できる。また、話者判別装置50に含まれる機能部の一部を別の集積回路や電子回路とすることもできる。例えば、集積回路としては、ASIC(Application Specific Integrated Circuit)が挙げられる。また、電子回路としては、CPU(Central Processing Unit)やMPU(Micro Processing Unit)などが挙げられる。

【0063】

続いて、本実施例に係る会話分析装置の処理の流れについて説明する。なお、ここでは、会話分析装置10によって実行される(1)会話分析処理を説明した後に、話者判別装置50によって実行される(2)有効化処理を説明する。

【0064】

(1) 会話分析処理

図9及び図10は、実施例1に係る会話分析処理の手順を示すフローチャートである。この会話分析処理は、一例として、図示しない入力部から会話分析処理を実行する指示を受け付けた場合に処理が起動する。

【0065】

図9に示すように、取得部51は、第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cを取得する(ステップS101)。そして、フレーム化部52は、第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cそれぞれの長さが同一であるか否かを判定する(ステップS102)。なお、ここで言う「同一」は、長さの差が許容誤差範囲内である場合も含む。

【0066】

このとき、各々の音声データの長さが同一でない場合(ステップS102否定)には、フレーム化部52は、エラーメッセージを図示しない表示部に出力し(ステップS103)、処理を終了する。

【0067】

一方、各々の音声データの長さが同一である場合(ステップS102肯定)には、フレーム化部52は、第1の音声データ12A、第2の音声データ12B及び第3の音声データ12Cをフレーム化する(ステップS104)。

【0068】

その後、第1の識別部53は、自己相関係数のピークの数、自己相関係数のピークの最

10

20

30

40

50

大値およびスペクトルエントロピーの3つの特徴量を各々の音声データごとに抽出する(ステップS105)。そして、第1の識別部53は、各々の音声データごとに抽出した3つの特徴量それぞれの平均値および標準偏差を算出する(ステップS106)。

【0069】

続いて、第1の識別部53は、変数Nに0を設定し(ステップS107)、隠れマルコフモデルにおける有声音および無声音の状態遷移について初期の状態遷移確率 P_i を設定する(ステップS108)。

【0070】

そして、第1の識別部53は、変数Nの値を1つインクリメントする(ステップS109)。このとき、変数Nの値が5以上でない場合(ステップS110否定)には、第1の識別部53は、各々の音声データごとに抽出した上記の3つの特徴量、並びに、各特徴量の平均値および標準偏差を観測結果とし、EM法を用いて、状態遷移確率 P_i を算出し(ステップS111)、ステップS109へ移行する。

【0071】

一方、変数Nの値が5以上である場合(ステップS110肯定)には、第1の識別部53は、各々の音声データごとに抽出した上記の3つの特徴量、並びに、各特徴量の平均値および標準偏差を観測結果とし、EM法を用いて、状態遷移確率 P_i を算出する(ステップS112)。

【0072】

そして、第1の識別部53は、各々の音声データごとに抽出した上記の3つの特徴量、並びに、各特徴量の平均値および標準偏差を観測結果とし、ビタビアルゴリズムを用いて、観測確率 P_o を算出する(ステップS113)。

【0073】

その後、第1の識別部53は、各々の音声データごとに抽出した上記の3つの特徴量に基づいて、ビタビアルゴリズムを用いて、次のような処理を行う。すなわち、第1の識別部53は、発話が行われている各フレームにおいて、発話されている音が有声音であるか、あるいは無声音であるかを識別する。そして、第1の識別部53は、有声音が検出された領域を有声音領域とし、無声音が検出された領域を無声音領域とする(ステップS114)。

【0074】

ここで、有効化部54は、各音声データにおける同一区間のフレームで有声音領域が重複する場合に、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレームの識別結果を有効化する「有効化処理」を実行する(ステップS115)。

【0075】

その後、第2の識別部55は、有効化部54による有効化後の有声音および無声音に基づいて、ビタビアルゴリズムを用いて、沈黙状態であるか、あるいは発話状態であるかを検出することで、沈黙領域および発話領域を識別する(ステップS116)。

【0076】

続いて、抽出部13は、図10に示すように、ある話者が発話したと特定されたフレームから、有声音領域の数、有声音領域の長さの平均値および有声音領域の長さの標準偏差を算出する(ステップS117)。

【0077】

さらに、抽出部13は、ある話者が発話したと特定されたフレームから、発話領域の数、発話領域の長さの平均値および発話領域の長さの標準偏差を算出する(ステップS118)。その後、抽出部13は、ある話者の沈黙領域のフレームから、沈黙領域の数、沈黙領域の長さの平均値および沈黙領域の長さの標準偏差を算出する(ステップS119)。

【0078】

そして、抽出部13は、会話全体の時間の長さに対するある話者の発話時間の長さの割合を算出する(ステップS120)。さらに、抽出部13は、他の話者の発話時間に対す

10

20

30

40

50

るある話者の発話時間の割合を算出する（ステップS121）。

【0079】

続いて、抽出部13は、ある話者が発話したと特定されたフレームから、音量の標準偏差およびスペクトルエントロピーの標準偏差を算出する（ステップS122）。抽出部13は、ある話者が発話したと特定されたフレームから算出した音量の標準偏差と、スペクトルエントロピーの標準偏差との和を、変化の度合いとして算出する（ステップS123）。

【0080】

そして、全ての話者の会話特性を抽出するまで（ステップS124否定）、上記のステップS117～ステップS123までの処理を繰り返し実行する。その後、全ての話者の会話特性を抽出すると（ステップS124肯定）、分析部14は、抽出部13によって抽出された会話特性に基づいて、会話スタイルを分析する（ステップS125）。最後に、分析部14は、分析結果を所定の出力先の装置へ出力し（ステップS126）、処理を終了する。

【0081】

（2）有効化処理

図11は、実施例1に係る有効化処理の手順を示すフローチャートである。この有効化処理は、図9に示したステップS115に対応する処理であり、有声音領域および無声音領域が識別された後に処理が起動する。

【0082】

図11に示すように、有効化部54は、各音声データにおける同一区間のフレームで第1の識別部53による識別結果を比較する（ステップS301）。このとき、同一区間のフレームで有声音領域が重複する場合（ステップS302肯定）には、有効化部54は、当該有声音領域と識別されたフレームのエネルギーを演算する（ステップS303）。なお、同一区間のフレームで有声音領域が重複しない場合（ステップS302否定）には、ステップS306へ移行する。

【0083】

そして、有効化部54は、当該同一区間で有声音領域と識別されたフレームのうち最大エネルギーを持つフレームを特定する（ステップS304）。その上で、有効化部54は、最大エネルギーを持つフレーム以外の識別結果を有声音領域から無声音領域に置き換える（ステップS305）。

【0084】

その後、各音声データ間で同一区間のフレームを全て処理するまで（ステップS306否定）、上記のステップS301～ステップS305までの処理を繰り返し実行する。そして、各音声データ間で同一区間のフレームを全て処理すると（ステップS306肯定）、処理を終了する。

【0085】

[実施例1の効果]

上述してきたように、本実施例に係る話者判別装置50は、各音声データにおける同一区間のフレームで有声音領域が重複して識別された場合に、最大のエネルギーを持つフレームの識別結果だけを有効化して各々の音声データの発話領域および沈黙領域を識別する。このため、本実施例に係る話者判別装置50は、各音声データを構成する同一区間のフレーム間で閾値を用いて判定せずとも、話者を判別することができる。さらに、本実施例に係る話者判別装置50では、話者の判別に複雑なアルゴリズムを用いる必要はなく、事前に学習を行う必要もない。したがって、本実施例に係る話者判別装置50によれば、話者の判別を簡易かつ正確に行うことができる。

【0086】

また、本実施例に係る話者判別装置50は、各音声データにおける同一区間のフレームで有声音領域が単独で識別された場合には、エネルギーの大小に関係なく、有声音領域と識別された識別結果を維持する。一般に、発話は、有声音と無声音が混在して構成される

10

20

30

40

50

ので、複数の話者によって同時に発話された場合でも、同時発話で有声音領域が完全に重複する可能性は低く、有声音領域が単独で識別される機会が残る可能性は高い。それゆえ、本実施例に係る話者判別装置 50 では、話者が発話する音量に開きがある場合でも、同時発話を判別することもできる。

【0087】

さらに、本実施例に係る話者判別装置 50 は、当該同一区間で有声音領域と識別されたフレームのうち最大のエネルギーを持つフレーム以外の識別結果を無声音領域に置き換える。このため、本実施例に係る話者判別装置 50 では、識別情報の置換という簡易な処理によって最大のエネルギーを持つフレームの識別結果だけを有効化できる結果、話者の判別を簡易に実現できる。

【実施例 2】

【0088】

さて、これまで開示の装置に関する実施例について説明したが、本発明は上述した実施例以外にも、種々の異なる形態にて実施されてよいものである。そこで、以下では、本発明に含まれる他の実施例を説明する。

【0089】

[エネルギー]

例えば、上記の実施例 1 では、最大エネルギーを持つフレームの識別結果だけを有効化する場合を例示したが、エネルギーに関連する他の指標が最大となるフレームの識別結果だけを有効化することもできる。一例としては、開示の装置は、フレームで観測される振幅の最大値および最小値の差が最大であるフレームの識別結果だけを有効化することもできる。この場合には、エネルギーの演算処理よりも簡易な演算により、識別結果の置換を実現できる。

【0090】

[マイク]

また、上記の実施例 1 では、接話型マイクロホンを適用する場合を例示したが、開示の装置はこれに限定されず、必ずしもマイクを装着する話者以外の他の話者をマイクから遠ざける必要はない。例えば、指向性を持つマイクを適用することができる。この場合には、話者 A が発話する方向の感度が他の方向の感度よりも強くなるように話者 A または指向性マイクを配置し、また、話者 B および話者 C についても同様にして指向性マイクを用い

ればよい。なお、指向性マイクを用いる場合についても、話者は複数であればよく、2 人であっても 4 人以上であっても開示の装置を適用できる。

【0091】

[分散および統合]

また、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。例えば、話者判別装置 50、抽出部 13 または分析部 14 を会話分析装置の外部装置としてネットワーク経由で接続するようにしてもよい。また、話者判別装置 50、抽出部 13 または分析部 14 を別の装置がそれぞれ有し、ネットワーク接続されて協働することで、上記の話者判別装置の機能を実現するよう

にしてもよい。

【0092】

[話者判別プログラム]

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。そこで、以下では、図 12 を用いて、上記の実施例と同様の機能を有する話者判別プログラムを実行するコンピュータの一例について説明する。

【0093】

図 12 は、実施例 1 及び実施例 2 に係る話者判別プログラムを実行するコンピュータの

10

20

30

40

50

一例について説明するための図である。図12に示すように、コンピュータ100は、操作部110aと、スピーカ110bと、マイク110cと、ディスプレイ120と、通信部130とを有する。さらに、このコンピュータ100は、CPU150と、ROM160と、HDD170と、RAM180とを有する。これら110~180の各部はバス140を介して接続される。

【0094】

HDD170には、図12に示すように、上記の実施例1で示した取得部51と、フレーム化部52と、第1の識別部53と、有効化部54と、第2の識別部55と同様の機能を発揮する話者判別プログラム170aが予め記憶される。この話者判別プログラム170aについては、図1に示した各々の取得部51、フレーム化部52、第1の識別部53、有効化部54及び第2の識別部55の各構成要素と同様、適宜統合又は分離しても良い。すなわち、HDD170に格納される各データは、常に全てのデータがHDD170に格納される必要はなく、処理に必要なデータのみがHDD170に格納されれば良い。

10

【0095】

そして、CPU150が、話者判別プログラム170aをHDD170から読み出してRAM180に展開する。これによって、図12に示すように、話者判別プログラム170aは、話者判別プロセス180aとして機能する。この話者判別プロセス180aは、HDD170から読み出した各種データを適宜RAM180上の自身に割り当てられた領域に展開し、この展開した各種データに基づいて各種処理を実行する。なお、話者判別プロセス180aは、図1に示した取得部51、フレーム化部52、第1の識別部53、有効化部54及び第2の識別部55にて実行される処理、例えば図9~図11に示す処理を含む。また、CPU150上で仮想的に実現される各処理部は、常に全ての処理部がCPU150上で動作する必要はなく、処理に必要な処理部のみが仮想的に実現されれば良い。

20

【0096】

なお、上記の話者判別プログラム170aについては、必ずしも最初からHDD170やROM160に記憶させておく必要はない。例えば、コンピュータ100に挿入されるフレキシブルディスク、いわゆるFD、CD-ROM、DVDディスク、光磁気ディスク、ICカードなどの「可搬用の物理媒体」に各プログラムを記憶させる。そして、コンピュータ100がこれらの可搬用の物理媒体から各プログラムを取得して実行するようにしてもよい。また、公衆回線、インターネット、LAN、WANなどを介してコンピュータ100に接続される他のコンピュータまたはサーバ装置などに各プログラムを記憶させておき、コンピュータ100がこれらから各プログラムを取得して実行するようにしてもよい。

30

【符号の説明】

【0097】

- 10 会話分析装置
- 11 音声記憶部
- 12A 第1の音声データ
- 12B 第2の音声データ
- 12C 第3の音声データ
- 30A, 30B, 30C 接話マイク
- 31 登録部
- 50 話者判別装置
- 51 取得部
- 52 フレーム化部
- 53 第1の識別部
- 54 有効化部
- 55 第2の識別部

40

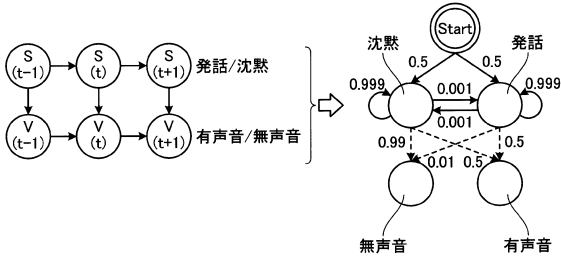
【図7】

図6に示した識別結果の置換結果の一例を示す図

時刻	第1音声データ	第2音声データ	第3音声データ
	V or U	V or U	V or U
12:00:00.000	V	U	U
12:00:00.010	V	U	U
12:00:00.020	U	V	U
12:00:00.030	U	U	U
⋮	⋮	⋮	⋮

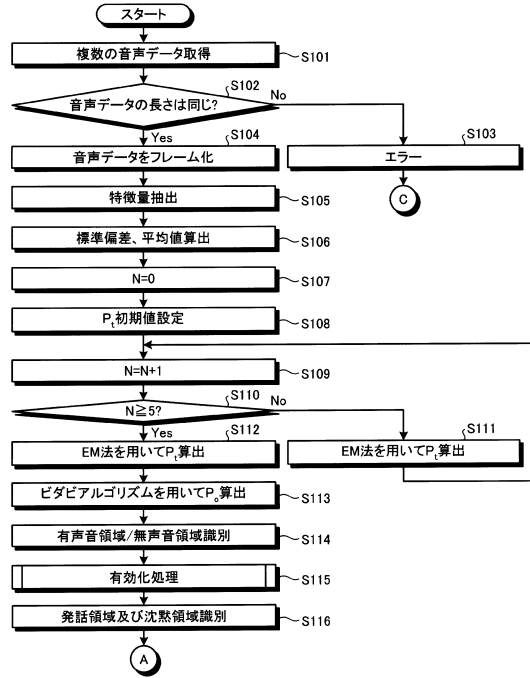
【図8】

隠れマルコフモデルにおける状態遷移図の一例を示す図



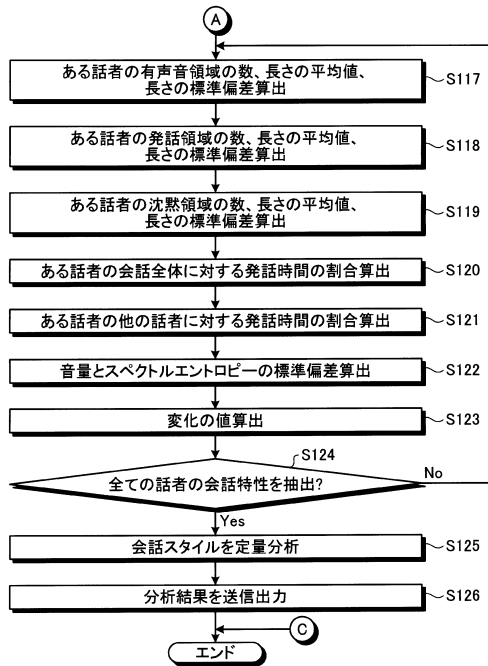
【図9】

実施例1に係る会話分析処理の手順を示すフローチャート



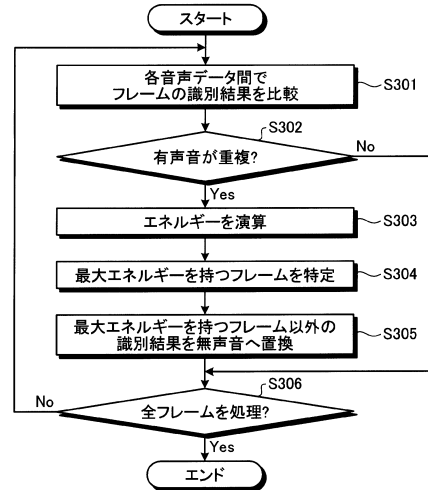
【図10】

実施例1に係る会話分析処理の手順を示すフローチャート

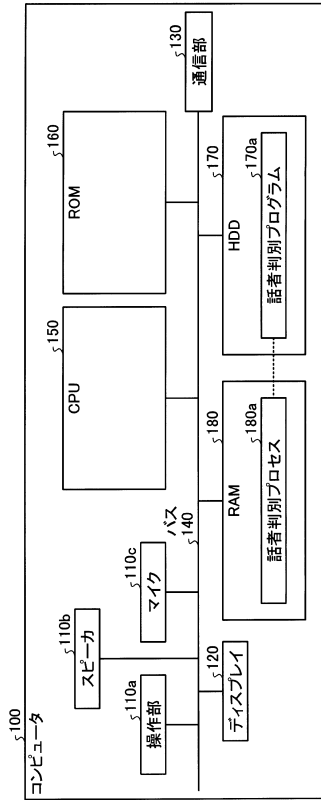


【図11】

実施例1に係る有効化処理の手順を示すフローチャート



実施例1及び実施例2に係る話者判別プログラムを実行するコンピュータの一例について説明するための図



フロントページの続き

(56)参考文献 特開2006-208482(JP,A)
特開2001-343983(JP,A)
特開平07-056598(JP,A)
特開2001-343985(JP,A)
特開昭58-095399(JP,A)
特開2006-086877(JP,A)
米国特許第05197113(US,A)

(58)調査した分野(Int.Cl., DB名)

G10L 15/00 - 15/34
G10L 25/00 - 25/93