(19) **日本国特許庁(JP)**

(12) 特許公報(B2)

(11)特許番号

特許第4811993号 (P4811993)

(45) 発行日 平成23年11月9日(2011.11.9)

(24) 登録日 平成23年9月2日(2011.9.2)

(51) Int.Cl.			FΙ		
G10L	11/00	(2006.01)	G1OL	11/00	402A
G10L	21/04	(2006.01)	G10L	11/00	1 O 1 D
G10L	15/20	(2006.01)	G10L	21/04	120B
			G10L	15/20	360Z

請求項の数 10 (全 54 頁)

特願2005-241264 (P2005-241264) (21) 出願番号 (22) 出願日 平成17年8月23日 (2005.8.23) (65) 公開番号 特開2007-57692 (P2007-57692A) (43) 公開日 平成19年3月8日 (2007.3.8) 平成20年3月26日 (2008.3.26) 審査請求日

(出願人による申告) 平成17年度独立行政法人情報通 信研究機構、研究テーマ「人間情報コミュニケーション の研究開発」に関する委託研究、産業活力再生特別措置 ||(72)発明者 渡辺 秀行 法第30条の適用を受ける特許出願

||(73)特許権者 393031586

株式会社国際電気通信基礎技術研究所 京都府相楽郡精華町光台二丁目2番地2

||(74)代理人 100115749

弁理士 谷川 英和

(72) 発明者 山田 玲子

> 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

(72) 発明者 田川 博章

> 京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

> > 最終頁に続く

(54) 【発明の名称】音声処理装置、およびプログラム

(57)【特許請求の範囲】

【請求項1】

比較される対象の音声に関するデータであり、1以上の音韻毎のデータであり、状態を識 別する2以上の状態識別子と状態間の遷移確率の情報を有する教師データを1以上格納し ている教師データ格納部と、

音声を受け付ける音声受付部と、

第一サンプリング周波数を格納している第一サンプリング周波数格納部と、

前記教師データのフォルマント周波数である教師データフォルマント周波数を格納してい る教師データフォルマント周波数格納部と、

前記音声受付部が受け付けた音声の話者である評価対象者のフォルマント周波数である評 価対象者フォルマント周波数を格納している評価対象者フォルマント周波数格納部と、

第二サンプリング周波数「前記第一サンプリング周波数/(教師データフォルマント周波 数/評価対象者フォルマント周波数)」で、前記音声受付部が受け付けた音声に対して、 サンプリング処理を行い、第二音声データを得る声道長正規化処理部と、

前記第二音声データを処理する音声処理部を具備し、

前記音声処理部は、

前記第二音声データを、フレームに区分するフレーム区分手段と、

前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム 音声データ取得手段と、

前記教師データと前記1以上のフレーム音声データに基づいて、前記音声受付部が受け付

けた音声の評定を行う評定手段と、

前記評定手段における評定結果を出力する出力手段を具備し、

前記評定手段は、

<u>前記1以上のフレーム音声データのうちの少なくとも一のフレーム音声データに対する最</u> 適状態を決定する最適状態決定手段と、

前記最適状態決定手段が決定した最適状態における確率値と、当該確率値に対応するフレームの全状態における確率値の総和との比率を用いて、事後確率を示す確率値である最適 状態確率値を取得する最適状態確率値取得手段と、

前記最適状態確率値取得手段が取得した最適状態確率値をパラメータとして音声の評定値 を算出する評定値算出手段を具備する音声処理装置。

【請求項2】

比較される対象の音声に関するデータであり、1以上の音韻毎のデータであり、状態を識別する2以上の状態識別子と状態間の遷移確率の情報を有する教師データを1以上格納している教師データ格納部と、

音声を受け付ける音声受付部と、

第一サンプリング周波数を格納している第一サンプリング周波数格納部と、

<u>前記教師データのフォルマント周波数である教師データフォルマント周波数を格納してい</u>る教師データフォルマント周波数格納部と、

前記音声受付部が受け付けた音声の話者である評価対象者のフォルマント周波数である評価対象者フォルマント周波数を格納している評価対象者フォルマント周波数格納部と、

第二サンプリング周波数「前記第一サンプリング周波数 / (教師データフォルマント周波数 / 評価対象者フォルマント周波数)」で、前記音声受付部が受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理部と、

前記第二音声データを処理する音声処理部を具備し、

前記音声処理部は、

前記第二音声データを、フレームに区分するフレーム区分手段と、

<u>前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム</u> 音声データ取得手段と、

<u>前記教師データと前記1以上のフレーム音声データに基づいて、前記音声受付部が受け付</u>けた音声の評定を行う評定手段と、

前記評定手段における評定結果を出力する出力手段を具備し、

前記評定手段は、

前記1以上のフレーム音声データの最適状態を決定する最適状態決定手段と、

前記最適状態決定手段が決定した各フレームの最適状態を有する音韻全体の状態における 2 以上の事後確率を示す確率値を、発音区間毎に取得する発音区間フレーム音韻確率値取 得手段と、

前記発音区間フレーム音韻確率値取得手段が取得した2以上の確率値の総和を、フレーム毎に算出し、当該フレーム毎の確率値の総和に基づいて、発音区間毎の確率値の総和の時間平均値を1以上算出し、当該1以上の時間平均値を用いて音声の評定値を算出する評定値算出手段を具備する音声処理装置。

【請求項3】

<u>前記第一サンプリング周波数で、前記音声受付部が受け付けた音声をサンプリングし、第</u> 一音声データを取得するサンプリング部をさらに具備し、

前記声道長正規化処理部は、

第二サンプリング周波数「前記第一サンプリング周波数 / (教師データフォルマント周波数 / 評価対象者フォルマント周波数) 」で、前記第一音声データに対して、リサンプリング処理を行い、第二音声データを得る請求項 1 または請求項 2 記載の音声処理装置。

【請求項4】

前記音声処理部は、

前記フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する特

20

10

30

殊音声検知手段をさらに具備し、

前記評定手段は、

前記教師データと前記入力音声データと前記特殊音声検知手段における検知結果に基づいて、前記音声受付部が受け付けた音声の評定を行う請求項<u>1から請求項3いずれか</u>記載の音声処理装置。

【請求項5】

前記特殊音声検知手段は、

無音を示すHMMに基づくデータである無音データを格納している無音データ格納手段と

前記入力音声データおよび前記無音データに基づいて、無音の区間を検出する無音区間検 出手段を具備<u>し、</u>

前記評定手段は、

<u>前記無音データの区間のデータを用いずに、音声の評定値を算出</u>する請求項<u>4</u>記載の音声 処理装置。

【請求項6】

前記特殊音声検知手段は、

一の音素の少なくとも最終フレームを含む後半部および当該音素の次の音素の少なくとも 第一フレームを含む前半部の評定値が所定値より低い場合、または一の音素の少なくとも 最終フレームを含む後半部および当該音素の次の音素の少なくとも第一フレームを含む前 半部の評定値が所定値より低くかつ無音が挿入されていないと判断した場合、または、一 の音素の少なくとも最終フレームを含む後半部および当該音素の次の音素の少なくとも第 一フレームを含む前半部の評定値が所定値より低くかつ他の音韻 H M M に対する確率値を 算出し、所定の値より高い確率値を得た音韻を検出した場合、または評定値が所定値より 低い音素が連続した場合に、音素の挿入があったと判断し、

前記評定手段は、

前記特殊音声検知手段が<u>音素の挿入があったと判断</u>した場合に、少なくとも音素の挿入が あった旨を示す評定結果を構成する請求項 4 記載の音声処理装置。

【請求項7】

前記特殊音声検知手段は、

一の音素の評定値が所定の値より低くかつ当該音素の直前の音素および当該音素の直後の 音素の評定値が所定の値より高い場合、または一の音素の評定値が所定の値より低く、か つ想定していない音素のHMMに基づいて算出された評定値が所定の値より高い場合に、 音韻の置換があったと判断し、

または、

一の音素の評定値が所定の値より低くかつ当該音素の直前の音素および当該音素の直後の音素の評定値が所定の値より高い場合、または一の音素の評定値が所定の値より低くかつ当該音素の直前の音素および当該音素の直後の音素の評定値が所定の値より高くかつ当該音素の区間長が所定の長さよりも短い場合、または直前の音素に対応する確率値、または直後の音素に対応する確率値が、当該一の音素の確率値より高い場合に、音韻の欠落があったと判断し、

前記評定手段は、

前記特殊音声検知手段が<u>音素の置換があったと判断</u>した場合に、少なくとも音素の置換が あった旨を示す評定結果を構成し、

または、

前記特殊音声検知手段が音素の欠落があったと判断した場合に、少なくとも音素の欠落が あった旨を示す評定結果を構成する請求項 4 記載の音声処理装置。

【請求項8】

前記音声処理装置は、カラオケ評価装置であって、

前記音声受付部は、

評価対象者の歌声の入力を受け付け、

20

10

30

40

前記音声処理部は、

前記歌声を評価する請求項1から請求項7いずれか記載の音声処理装置。

【請求項9】

第一サンプリング周波数を格納している第一サンプリング周波数格納部と、

前記教師データのフォルマント周波数である教師データフォルマント周波数を格納してい る教師データフォルマント周波数格納部と、

前記音声の話者である評価対象者のフォルマント周波数である評価対象者フォルマント周 波数を格納している評価対象者フォルマント周波数格納部と、

第二サンプリング周波数「前記第一サンプリング周波数/(教師データフォルマント周波 数/評価対象者フォルマント周波数)」で、前記受け付けた音声に対して、サンプリング 処理を行い、第二音声データを得る声道長正規化処理部を具備し、

10

前記音声処理部は、

前記第二音声データを、フレームに区分するフレーム区分手段と、

前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム 音声データ取得手段と、

前記教師データと前記1以上のフレーム音声データに基づいて、前記音声受付部が受け付 けた音声の評定を行う評定手段と、

前記評定手段における評定結果を出力する出力手段を具備し、

前記評定手段は、

前記1以上のフレーム音声データのうちの少なくとも一のフレーム音声データに対する最 適状態を決定する最適状態決定手段と、

前記最適状態決定手段が決定した最適状態における確率値と、当該確率値に対応するフレ ームの全状態における確率値の総和との比率を用いて、事後確率を示す確率値である最適 状態確率値を取得する最適状態確率値取得手段と、

前記最適状態確率値取得手段が取得した最適状態確率値をパラメータとして音声の評定値 を算出する評定値算出手段を具備するデジタルシグナルプロセッサ。

【 請 求 項 1 0 】

第一サンプリング周波数を格納している第一サンプリング周波数格納部と、

前記教師データのフォルマント周波数である教師データフォルマント周波数を格納してい る教師データフォルマント周波数格納部と、

前記音声の話者である評価対象者のフォルマント周波数である評価対象者フォルマント周 波数を格納している評価対象者フォルマント周波数格納部と、

第ニサンプリング周波数「前記第一サンプリング周波数/(教師データフォルマント周波 数/評価対象者フォルマント周波数)」で、前記受け付けた音声に対して、サンプリング 処理を行い、第二音声データを得る声道長正規化処理部を具備し、

前記音声処理部は、

前記第二音声データを、フレームに区分するフレーム区分手段と、

前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム 音声データ取得手段と、

前記教師データと前記1以上のフレーム音声データに基づいて、前記音声受付部が受け付 けた音声の評定を行う評定手段と、

前記評定手段における評定結果を出力する出力手段を具備し、

前記評定手段は、

前記1以上のフレーム音声データの最適状態を決定する最適状態決定手段と、

前記最適状態決定手段が決定した各フレームの最適状態を有する音韻全体の状態における 2 以上の事後確率を示す確率値を、発音区間毎に取得する発音区間フレーム音韻確率値取 得手段と、

前記発音区間フレーム音韻確率値取得手段が取得した2以上の確率値の総和を、フレーム 毎に算出し、当該フレーム毎の確率値の総和に基づいて、発音区間毎の確率値の総和の時 間平均値を1以上算出し、当該1以上の時間平均値を用いて音声の評定値を算出する評定

20

30

40

値算出手段を具備するデジタルシグナルプロセッサ。

【発明の詳細な説明】

【技術分野】

[00001]

本発明は、入力された音声を評価したり、入力された音声を認識したりする音声処理装置等に関するものである。

【背景技術】

[0002]

従来の技術として、以下の音声処理装置がある(特許文献 1 参照)。本音声処理装置は、語学学習装置であり、当該語学学習装置は、学習者が選択した役割の発音をレファランスデータと比較して一致度によって点数化して表示し、点数によって適当な次の画面を自動に表示することにより、学習能率を向上させる装置である。本従来の音声処理装置は、入力された音声信号は音声認識技術により分析された後、学習者発音のスペクトルと抑揚とが学習者発音表示ボックスに表れるという構成になっている。そして、従来の技術においては、標準音データと学習者の発音のスペクトル、および抑揚が比較されて点数が表示される。

[0003]

また、従来の技術として、以下の音声処理装置がある(特許文献 2 参照)。本音声処理 装置は歌唱音声評価装置であり、本歌唱音声評価装置は、歌唱音声の周波数成分を抽出す る抽出手段と、当該抽出された周波数成分から基本周波数成分と倍音周波数成分とをそれ ぞれ抽出する特定周波数成分抽出手段と、特定周波数成分抽出手段によって抽出された基 本周波数成分に対する倍音周波数成分の比率に応じて、歌唱音声の評価を示す評価値を算 出する評価手段とを備える。そして、本歌唱音声評価装置は、歌唱音声の周波数成分に基 づいてその声質の良否を適正に評価し、これを歌唱音声の採点結果に反映させることによ り、歌唱音声の採点をより人間の感性に近づけることを狙いとしている。

[0004]

さらに、従来の技術として、以下の音声処理装置がある(特許文献3参照)。本音声処理装置は音声認識装置であり、入力音声パターンと標準パターンを、DP法を用いて照合し、最も照合距離の小さい標準パターンを認識結果とする音声認識装置であり、照合結果を用いて入力パターンを音素に分割し、各音素の継続時間と標準継続時間とのずれの分散を計算し、これを照合距離に付加することで距離を補正することを特徴とする。そして、分割部で照合結果を用いて音素に分割し、時間長ずれ計算部で標準継続時間とのずれの分散を計算し、距離補正部で照合距離を補正するように構成する。また、本音声認識装置は、時間長のずれを計算する対象音素を選択する音素選択部、距離補正する対象単語を選択する単語選択部を有し、単語の認識性能を高できる、というものである。

【特許文献1】特開2003-228279(第1頁、第1図等)

【特許文献2】特開2005-107088(第1頁、第1図等)

【特許文献3】特開平6-4096(第1頁、第1図等)

【発明の開示】

【発明が解決しようとする課題】

[0005]

しかしながら、特許文献 1 や特許文献 2 の従来の技術においては、音声(歌声も含む)の話者である評価対象者の話者特性に応じた音声処理が行えず、その結果、精度の高い音声処理ができなかった。具体的には、従来の技術においては、例えば、評価対象者の声道長の違いにより、スペクトル包絡が高周波数域または低周波数域に伸縮するが、従来の発音評定装置や歌唱音声評価装置などの音声処理装置において、かかるスペクトル包絡の伸縮により、評価結果が異なる。つまり、従来の技術においては、同様の上手さの発音や歌唱でも、評価対象者の声道長の違いにより、発音や歌唱の評価結果が異なり、精度の高い評価ができなかった。

[0006]

40

10

20

30

また、特許文献1の音声処理装置において、標準音データと学習者の発音のスペクトル、および抑揚が比較されて点数が表示される構成であるので、両者の類似度の評定の精度が低く、また、リアルタイムに高速に点数を表示するためには、処理能力が極めて高いCPU、多量のメモリが必要であった。

[0007]

また、特許文献1の音声処理装置において、無音区間があれば、類似度が低く評価されると考えられ、評価の精度が低かった。また、音素の置換や挿入や欠落など、特殊な事象が発生していることを検知できなかった。

[0008]

さらに、例えば、特許文献3に示すような音声認識処理を行う音声処理装置において、評価対象者の声道長の違いにより、スペクトル包絡の伸縮が生じるが、かかる評価対象者の話者特性に応じた音声認識処理を行っておらず、精度の高い音声認識ができなかった。

【課題を解決するための手段】

[0009]

本第一の発明の音声処理装置は、比較される対象の音声に関するデータであり、1以上の音韻毎のデータである教師データを1以上格納している教師データ格納部と、音声を受け付ける音声受付部と、第一サンプリング周波数を格納している第一サンプリング周波数 格納部と、前記第一サンプリング周波数で、前記音声受付部が受け付けた音声をサンプリング部と、前記教師データのフォルマント周波数である教師データフォルマント周波数を格納している教師データフォルマント周波数格納部と、前記音声受付部が受け付けた音声の話者である評価対象者のフォルマント周波数のある評価対象者フォルマント周波数を格納している評価対象者フォルマント周波数が割さいる評価対象者フォルマント周波数が割さいる評価対象者フォルマント周波数が割さいる評価対象者フォルマント周波数が割さいる評価対象者フォルマント周波数が割が受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理部と、前記第二音声データを処理する音声処理部を具備する音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声処理ができる。

[0010]

また、本第二の発明の音声処理装置は、第一の発明に対して、前記音声処理部は、前記第二音声データを、フレームに区分するフレーム区分手段と、前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム音声データ取得手段と、前記教師データと前記1以上のフレーム音声データに基づいて、前記音声受付部が受け付けた音声の評定を行う評定手段と、前記評定手段における評定結果を出力する出力手段を具備する音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声の評定ができる。

[0011]

また、本第三の発明の音声処理装置は、第二の発明に対して、前記評定手段は、前記1以上のフレーム音声データのうちの少なくとも一のフレーム音声データに対する最適状態を決定する最適状態決定手段と、前記最適状態決定手段が決定した最適状態における確率値を取得する最適状態確率値取得手段と、前記最適状態確率値取得手段が取得した確率値をパラメータとして音声の評定値を算出する評定値算出手段を具備する音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声の評定ができる。

[0012]

また、本第四の発明の音声処理装置は、第二の発明に対して、前記評定手段は、前記1以上のフレーム音声データの最適状態を決定する最適状態決定手段と、前記最適状態決定手段が決定した各フレームの最適状態を有する音韻全体の状態における1以上の確率値を、発音区間毎に取得する発音区間フレーム音韻確率値取得手段と、前記発音区間フレーム音韻確率値取得手段が取得した1以上の発音区間毎の1以上の確率値をパラメータとして音声の評定値を算出する評定値算出手段を具備する音声処理装置である。

10

20

30

40

かかる構成により、評価対象者の話者特性に応じた、さらに精度の高い音声の評定ができる。

[0013]

また、本第五の発明の音声処理装置は、第二の発明に対して、前記音声処理部は、前記 フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する特殊音 声検知手段をさらに具備し、前記評定手段は、前記教師データと前記入力音声データと前 記特殊音声検知手段における検知結果に基づいて、前記音声受付部が受け付けた音声の評 定を行う音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声の評定ができ、かつ特殊音声を検知し、かかる特殊音声に対応した音声の評定ができる。

[0014]

また、本第六の発明の音声処理装置は、第五の発明に対して、前記特殊音声検知手段は、無音を示すHMMに基づくデータである無音データを格納している無音データ格納手段と、前記入力音声データおよび前記無音データに基づいて、無音の区間を検出する無音区間検出手段を具備する音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声の評定ができ、かつ 無音区間を検知し、かかる無音区間に対応した音声の評定ができる。

[0015]

また、本第七の発明の音声処理装置は、第五の発明に対して、前記特殊音声検知手段は、一の音素の後半部および当該音素の次の音素の前半部の評定値が所定の条件を満たすことを検知し、前記評定手段は、前記特殊音声検知手段が前記所定の条件を満たすことを検知した場合に、少なくとも音素の挿入があった旨を示す評定結果を構成する音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声の評定ができ、かつ音素の挿入を検知し、かかる音素の挿入に対応した音声の評定ができる。

[0016]

また、本第八の発明の音声処理装置は、第七の発明に対して、前記特殊音声検知手段は、一の音素の評定値が所定の条件を満たすことを検知し、前記評定手段は、前記特殊音声検知手段が前記所定の条件を満たすことを検知した場合に、少なくとも音素の置換または欠落があった旨を示す評定結果を構成する音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声の評定ができ、かつ音素の置換または欠落を検知し、かかる音素の置換または欠落に対応した音声の評定ができる。

[0017]

また、本第九の発明の音声処理装置は、第二から第八いずれかの発明に対して、前記音声処理装置は、カラオケ評価装置であって、前記音声受付部は、評価対象者の歌声の入力を受け付け、前記音声処理部は、前記歌声を評価する音声処理装置である。

かかる構成により、カラオケ評価装置として利用できる。

[0018]

また、本第十の発明の音声処理装置は、第九の発明に対して、前記フレーム区分手段は、前記音声をフレームに区分し、かつ、前記第二音声データをフレームに区分し、前記フレーム音声データ取得手段は、前記音声が区分されたフレーム毎の音声データである第一フレーム音声データを1以上得て、かつ前記第二音声データが区分されたフレーム毎の音声データである第二フレーム音声データを1以上得、前記評定手段は、前記教師データと前記1以上の第一フレーム音声データに基づいて、前記音声受付部が受け付けた音声の評定を行う第一評定手段と、前記教師データと前記1以上の第二フレーム音声データに基づいて、前記音声受付部が受け付けた音声の評定を行う第二評定手段と、前記第一評定手段における評定結果と前記第二評定手段における評定結果に基づいて、最終的な評定結果を得る評定結果取得手段とを具備する音声処理装置である。

かかる構成により、優れたカラオケ評価装置として利用できる。

10

30

20

50

[0019]

また、本第十一の発明の音声処理装置は、第九、第十いずれかの発明に対して、前記音声受付部は、所定の母音の音声を受け付けた後、評価対象者の歌声の入力を受け付け、前記サンプリング部は、前記第一サンプリング周波数で、前記母音の音声をもサンプリングし、前記サンプリングした母音の音声に基づいて、評価対象者のフォルマント周波数である評価対象者フォルマント周波数を取得する評価対象者フォルマント周波数取得部をさらに具備し、前記評価対象者フォルマント周波数格納部の評価対象者フォルマント周波数は、前記評価対象者フォルマント周波数取得部が取得した評価対象者フォルマント周波数である音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声の評定ができる。 また、本第十二の発明の音声処理装置は、第一の発明に対して、前記音声処理部は、前 記第二音声データに基づいて、音声認識処理を行う音声処理装置である。

かかる構成により、評価対象者の話者特性に応じた精度の高い音声認識ができる。

【発明の効果】

[0020]

本発明による音声処理装置によれば、評価対象者の話者特性に応じた精度の高い音声処理ができる。

【発明を実施するための最良の形態】

[0021]

以下、音声処理装置等の実施形態について図面を参照して説明する。なお、実施の形態において同じ符号を付した構成要素は同様の動作を行うので、再度の説明を省略する場合がある。

(実施の形態1)

[0022]

本実施の形態において、比較対象の音声と入力音声の類似度の評定を精度高く、かつ高速にできる音声処理装置について説明する。本音声処理装置は、音声(歌唱を含む)を評価する発音評定装置である。特に、本音声処理装置は、入力音声のフレームに対する最適状態の事後確率を、動的計画法を用いて算出することから、当該事後確率をDAP(Dynamic A Posteriori Probability)と呼び、DAPに基づく類似度計算法および発音評定装置をDAPSと呼ぶ。

[0023]

また、本実施の形態における音声処理装置は、例えば、語学学習や物真似練習やカラオケ評定などに利用できる。図1は、本実施の形態における音声処理装置のブロック図である。本音声処理装置は、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部104、第一サンプリング周波数格納部105、サンプリング部106、評価対象者フォルマント周波数取得部107、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部110を具備する。

音声処理部 1 1 0 は、フレーム区分手段 1 1 0 1、フレーム音声データ取得手段 1 1 0 2、評定手段 1 1 0 3、出力手段 1 1 0 4 を具備する。

評定手段1103は、最適状態決定手段11031、最適状態確率値取得手段1103 2、評定値算出手段11033を具備する。

[0024]

なお、音声処理装置は、キーボード 3 4 2 、マウス 3 4 3 などの入力手段からの入力を受け付ける。また、音声処理装置は、マイク 3 4 5 などの音声入力手段から音声入力を受け付ける。さらに、音声処理装置は、ディスプレイ 3 4 4 などの出力デバイスに情報を出力する。

[0025]

入力受付部 1 0 1 は、音声処理装置の動作開始を指示する動作開始指示や、入力した音声の評定結果の出力態様の変更を指示する出力態様変更指示や、処理を終了する終了指示などの入力を受け付ける。かかる指示等の入力手段は、テンキーやキーボードやマウスや

10

20

30

40

メニュー画面によるもの等、何でも良い。入力受付部 1 0 1 は、テンキーやキーボード等の入力手段のデバイスドライバーや、メニュー画面の制御ソフトウェア等で実現され得る

[0026]

[0027]

音声受付部103は、音声を受け付ける。音声受付部103は、例えば、マイク345のドライバーソフトで実現され得る。また、なお、音声受付部103は、マイク345とそのドライバーから実現されると考えても良い。音声は、マイク345から入力されても良いし、磁気テープやCD・ROMなどの記録媒体から読み出すことにより入力されても良い。

[0028]

教師データフォルマント周波数格納部104は、教師データのフォルマント周波数である教師データフォルマント周波数を格納している。教師データフォルマント周波数は、第一フォルマント周波数(F1)でも、第二フォルマント周波数(F2)でも、第三フォルマント周波数(F3)等でも良い。教師データフォルマント周波数格納部104の教師データフォルマント周波数は、予め格納されていても良いし、評価時に、動的に、教師データから取得しても良い。音声データからフォルマント周波数を取得する技術は、公知技術であるので説明を省略する。教師データフォルマント周波数格納部104は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。

[0029]

第一サンプリング周波数格納部 1 0 5 は、第一のサンプリング周波数である第一サンプリング周波数を格納している。第一サンプリング周波数は、評価対象者の音声を、最初にサンプリングする場合のサンプリング周波数である。第一サンプリング周波数格納部 1 0 5 は、不揮発性の記録媒体が好適であるが、揮発性の記録媒体でも実現可能である。

[0030]

サンプリング部 1 0 6 は、第一サンプリング周波数格納部 1 0 5 の第一サンプリング周波数で、音声受付部 1 0 3 が受け付けた音声をサンプリングし、第一音声データを取得する。なお、受け付けた音声をサンプリングする技術は公知技術であるので、詳細な説明を省略する。サンプリング部 1 0 6 は、通常、MPUやメモリ等から実現され得る。サンプリング部 1 0 6 の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0031]

評価対象者フォルマント周波数取得部107は、サンプリング部106が取得した第一音声データから、評価対象者のフォルマント周波数である評価対象者フォルマント周波数を取得する。評価対象者フォルマント周波数も、第一フォルマント周波数(F1)でも、

10

20

30

40

10

20

30

40

50

第二フォルマント周波数(F2)でも、第三フォルマント周波数(F3)でも良い。ただし、評価対象者フォルマント周波数と教師データフォルマント周波数は同一種のフォルマント周波数である。サンプリングして取得した第一音声データから、フォルマント周波数を取得する技術は公知技術であるので、詳細な説明を省略する。評価対象者フォルマント周波数取得部107は、通常、MPUやメモリ等から実現され得る。評価対象者フォルマント周波数取得部107の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0032]

評価対象者フォルマント周波数格納部108は、音声受付部103が受け付けた音声の話者である評価対象者のフォルマント周波数である評価対象者フォルマント周波数を、少なくとも一時的に格納している。評価対象者フォルマント周波数格納部108の評価対象者フォルマント周波数取得部107が取得したフォルマント周波数であるが、予め評価対象者フォルマント周波数を格納していても良い。評価対象者フォルマント周波数を格納していても良い。評価対象者フォルマント周波数格納部108に、予め評価対象者フォルマント周波数が格納されている場合、本音声処理装置において、評価対象者フォルマント周波数取得部107は不要である。評価対象者フォルマント周波数格納部108は、不揮発性の記録媒体でも、揮発性の記録媒体でも良い。

[0033]

声道長正規化処理部109は、第二サンプリング周波数で、音声受付部103が受け付 けた音声に対して、サンプリング処理を行い、第二音声データを得る。第二サンプリング 周波数は、「第一サンプリング周波数/(教師データフォルマント周波数/評価対象者フ ォルマント周波数)」で算出されるサンプリング周波数である。声道長正規化処理部 1 0 9は、音声受付部103が受け付けた音声をサンプリング処理して得られた第一音声デー タを、リサンプリング処理して第二音声データを得ることが好適であるが、音声受付部 1 0 3 が受け付けた音声をサンプリング処理し、直接的に第二音声データを得ても良い。直 接的に第二音声データを得る場合、例えば、サンプリング処理を行うハードウェアが可変 のサンプリング周波数でサンプリング処理を行えることが必要である。声道長正規化処理 部109は、通常、演算「教師データフォルマント周波数/評価対象者フォルマント周波 数」を行い、周波数スケール(「r」とする)を得る。そして、声道長正規化処理部10 9は、第一サンプリング周波数格納部105の第一サンプリング周波数(Fs)と「r」 に基づいて、演算「Fs/r」を行い、新しいサンプリング周波数(Fs/r)を得る。 この新しいサンプリング周波数(Fs/r)が第二サンプリング周波数である。次に、声 道長正規化処理部109は、第一音声データに対して、第二サンプリング周波数(Fs/ r)で、リサンプリング処理を行い、第二音声データを得る。声道長正規化処理部109 は、通常、MPUやメモリ等から実現され得る。声道長正規化処理部109の処理手順は 、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されて いる。但し、ハードウェア(専用回路)で実現しても良い。

[0034]

音声処理部110は、第二音声データを処理する。音声処理部110は、ここでは、評定処理である。ただし、音声処理部110は、音声認識や音声出力などの他の音声処理を行っても良い。音声出力は、単に、リサンプリング処理された音声を出力する処理である。なお、本実施の形態において、音声処理部110は、評定処理を行うものとして、説明する。音声処理部110は、通常、MPUやメモリ等から実現され得る。音声処理部110の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0035]

音声処理部110を構成しているフレーム区分手段1101は、第二音声データを、フレームに区分する。フレーム区分手段1101は、通常、MPUやメモリ等から実現され得る。フレーム区分手段1101の処理手順は、通常、ソフトウェアで実現され、当該ソ

フトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で 実現しても良い。

[0036]

音声処理部110を構成しているフレーム音声データ取得手段1102は、区分されたフレーム毎の音声データであるフレーム音声データを1以上得る。フレーム音声データ取得手段1102は、通常、MPUやメモリ等から実現され得る。フレーム音声データ取得手段1102の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0037]

音声処理部110を構成している評定手段1103は、教師データ格納部102の教師データと1以上のフレーム音声データに基づいて、音声受付部103が受け付けた音声の評定を行う。評定方法の具体例は、後述する。「音声受付部103が受け付けた音声を評定する」の概念には、第二音声データを評定することも含まれることは言うまでもない。評定手段1103は、通常、MPUやメモリ等から実現され得る。評定手段1103の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0038]

評定手段1103を構成している最適状態決定手段11031は、1以上のフレーム音声データのうちの少なくとも一のフレーム音声データに対する最適状態を決定する。最適状態決定手段11031は、例えば、全音韻HMMから、比較される対象(学習対象)の単語や文章などの音声を構成する1以上の音素に対応するHMMを取得し、当該取得した1以上のHMMから、音素の順序で連結したデータ(比較される対象の音声に関するデータであり、音韻毎の隠れマルコフモデルを連結したHMMに基づくデータ)を構成する。そして、構成した当該データ、および取得した特徴ベクトル系列を構成する各特徴ベクトルのよに基づいて、所定のフレームの最適状態(特徴ベクトルのよに対する最適状態)を決定する。なお、最適状態を決定するアルゴリズムは、例えば、Viterbiアルゴリズムである。また、教師データは、上述の比較される対象の音声に関するデータであり、音韻毎の隠れマルコフモデルを連結したHMMに基づくデータと考えても良いし、連結される前のデータであり、全音韻HMMのデータと考えても良い。

評定手段1103を構成している最適状態確率値取得手段11032は、最適状態決定手段11031が決定した最適状態における確率値を取得する。

[0039]

評定手段1103を構成している評定値算出手段11033は、最適状態確率値取得手段11032が取得した確率値をパラメータとして音声の評定値を算出する。評定値算出手段11033は、上記確率値を如何に利用して、評定値を算出するかは問わない。評定値算出手段11033は、例えば、最適状態確率値取得手段11032が取得した確率値と、当該確率値に対応するフレームの全状態における確率値の総和をパラメータとして音声の評定値を算出する。評定値算出手段11033は、ここでは、通常、フレームごとに評定値を算出する。

[0040]

最適状態決定手段11031、最適状態確率値取得手段11032、評定値算出手段11033は、通常、MPUやメモリ等から実現され得る。最適状態決定手段11031等の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0041]

出力手段1104は、評定手段1103における評定結果を出力する。出力手段110 4の出力態様は、種々考えられる。出力とは、ディスプレイへの表示、プリンタへの印字 、音出力、外部の装置への送信、記録媒体への蓄積等を含む概念である。出力手段110 4は、例えば、評定結果を視覚的に表示する。出力手段1104は、例えば、フレーム単 位、または/および音素・単語単位、または/および発声全体の評定結果を視覚的に表示 10

20

30

40

する。出力手段1104は、ディスプレイ344やスピーカー等の出力デバイスを含むと考えても含まないと考えても良い。出力手段1104は、出力デバイスのドライバーソフトと出力デバイス等で実現され得る。

次に、本音声処理装置の動作について図2、図3のフローチャートを用いて説明する。

[0042]

(ステップS201)入力受付部101は、音声処理装置の動作開始を指示する動作開始指示を受け付けたか否かを判断する。動作開始指示を受け付ければステップS202に行き、動作開始指示を受け付けなければステップS217に飛ぶ。

(ステップS202)音声受付部103は、音声を受け付けたか否かを判断する。音声を受け付ければステップS203に行き、音声を受け付けなければステップS216に飛ぶ。

10

[0043]

(ステップS203)サンプリング部106は、第一サンプリング周波数格納部105に格納されている第一サンプリング周波数を読み込み、当該第一サンプリング周波数で、音声受付部103が受け付けた音声をサンプリングし、第一音声データを得る。

[0044]

(ステップS204) 声道長正規化処理部109は、音声受付部103が受け付けた音声から、第二音声データを得る。かかる第二音声データを得る処理である声道長正規化処理の詳細については、図3のフローチャートを用いて、詳細に説明する。なお、声道長正規化処理は、個人差を吸収する評定のための前処理である。

20

(ステップS205)フレーム区分手段1101は、ステップS204で得た第二音声 データを図示しないバッファに一時格納する。

[0045]

(ステップS206)フレーム区分手段1101は、バッファに一時格納した第二音声データをフレームに区分する。かかる段階で、区分されたフレーム毎の音声データであるフレーム音声データが構成されている。フレーム区分手段1101が行うフレーム分割の処理は、例えば、フレーム音声データ取得手段1102がフレーム音声データを取り出す際の前処理であり、入力された音声のデータを、すべてのフレームに一度に分割するとは限らない。

30

(ステップS207)フレーム音声データ取得手段1102は、カウンタiに1を代入する。

[0046]

(ステップS208)フレーム音声データ取得手段1102は、 i 番目のフレームが存在するか否かを判断する。 i 番目のフレームが存在すればステップS209に行き、 i 番目のフレームが存在しなければステップS211に行く。

[0047]

(ステップS209)フレーム音声データ取得手段1102は、 i 番目のフレーム音声データを取得する。フレーム音声データの取得とは、例えば、当該分割された音声データを音声分析し、特徴ベクトルデータを抽出することである。なお、フレーム音声データは、例えば、入力された音声データをフレーム分割されたデータである。また、フレーム音声データは、例えば、当該分割された音声データから音声分析され、抽出された特徴ベクトルデータを有する。本特徴ベクトルデータは、例えば、三角型フィルタを用いたチャネル数24のフィルタバンク出力を離散コサイン変換したMFCCであり、その静的パラメータ、デルタパラメータおよびデルタデルタパラメータをそれぞれ12次元、さらに正規化されたパワーとデルタパワーおよびデルタデルタパワー(39次元)を有する。

40

(ステップS210)フレーム音声データ取得手段1102は、カウンタiを1、インクリメントする。ステップS208に戻る。

[0048]

(ステップS211)最適状態決定手段11031は、全フレームの最適状態を決定する。最適状態決定手段11031が最適状態を決定するアルゴリズムは、例えば、Vit

erbiアルゴリズムによる。Viterbiアルゴリズムは、公知のアルゴリズムであるので、詳細な説明は省略する。

[0049]

(ステップS212)最適状態確率値取得手段11032は、全フレームの全状態の前向き尤度、および後向き尤度を算出する。最適状態確率値取得手段11032は、例えば、全てのHMMを用いて、フォワード・バックワードアルゴリズムにより、前向き尤度、および後向き尤度を算出する。

(ステップS213)最適状態確率値取得手段11032は、ステップS212で取得した前向き尤度、および後向き尤度を用いて、最適状態の確率値(最適状態確率値)を、すべて算出する。

[0050]

(ステップS214)評定値算出手段11033は、ステップS213で算出した1以上の最適状態確率値から、1以上のフレームの音声の評定値を算出する。評定値算出手段11033は、例えば、1105最適状態確率値と、当該最適状態確率値に対応するフレームの全状態における確率値の総和をパラメータとして音声の評定値を算出する。詳細については、後述する。

[0051]

(ステップS215)出力手段1104は、ステップS214における評定結果(ここでは、音声の評定値)を、設定されている出力モードに従って、出力する。ステップS202に戻る。出力モードとは、評定値を数値で画面に表示するモード、評定値の遷移をグラフで画面に表示するモード、評定値を音声で出力するモード、評定値が所定の数値より低い場合に警告を示す情報を表示するモードなど、何でも良い。なお、ここでの出力モードは、ステップS218で設定されるモードである。

[0052]

(ステップ S 2 1 6) 音声受付部 1 0 3 は、タイムアウトか否かを判断する。つまり、音声受付部 1 0 3 は、所定の時間以上、音声の入力を受け付けなかったか否かを判断する。タイムアウトであればステップ S 2 0 1 に戻り、タイムアウトでなければステップ S 2 0 2 に戻る。

[0053]

(ステップS217)入力受付部101は、出力態様変更指示を受け付けたか否かを判断する。出力態様変更指示を受け付ければステップS218に行き、出力態様変更指示を受け付なければステップS219に飛ぶ。出力態様変更指示は、上述した出力モードを有する情報である。

(ステップS218)出力手段1104は、ステップS217で受け付けた出力態様変更指示が有する出力モードを示す情報を書き込み、出力モードを設定する。ステップS201に戻る。

(ステップS219)入力受付部101は、終了指示を受け付けたか否かを判断する。 終了指示を受け付ければ処理を終了し、終了指示を受け付なければステップS201に戻る。

なお、図 2 のフローチャートにおいて、本発音評定装置は、出力モードの設定機能を有 しなくても良い。

次に、ステップS204における声道長正規化処理の詳細について、図3のフローチャートを用いて説明する。

[0054]

(ステップS301)評価対象者フォルマント周波数取得部107は、サンプリング部106のサンプリング処理により得られた第一音声データから、評価対象者フォルマント周波数(Fi)を取得し、評価対象者フォルマント周波数格納部108に一時格納する。評価対象者フォルマント周波数は、例えば、第二フォルマント周波数(F2)である。

(ステップS302) 声道長正規化処理部109は、第一サンプリング周波数格納部105の第一サンプリング周波数(FS) を読み出す。

10

20

30

40

(ステップS303) 声道長正規化処理部109は、教師データフォルマント周波数格納部104の教師データフォルマント周波数を読み出す。

[0055]

(ステップS304) 声道長正規化処理部109は、ステップS301で取得した評価対象者フォルマント周波数と、ステップS303で読み出した教師データフォルマント周波数から周波数スケールを算出する。具体的には、声道長正規化処理部109は、演算「教師データフォルマント周波数/評価対象者フォルマント周波数」を行い、周波数スケール(r)を得る。

[0056]

(ステップS305) 声道長正規化処理部109は、ステップS302で読み出した第一サンプリング周波数(Fs) と周波数スケール(r) に基づいて、演算「Fs/r」を行い、第二サンプリング周波数(Fs/r)を得る。

[0057]

(ステップS306) 声道長正規化処理部109は、サンプリング部106がサンプリングして得た第一音声データに対して、第二サンプリング周波数(Fs/r)で、リサンプリング処理を行い、第二音声データを得る。なお、リサンプリング処理は公知技術であるので、詳細な説明を省略する。上位関数にリターンする。

[0058]

なお、図2、図3のフローチャートにおいて、声道長正規化処理を行う対象の音声と、評価対象の音声が異なっても良い。つまり、例えば、音声受付部103は、所定の1以上の母音(例えば、「う」)の音声を受け付けた後、評価対象者の音声を受け付け、評価対象者フォルマント周波数取得部107は、当該1以上の母音の音声に基づいて、評価対象者フォルマント周波数を取得し、声道長正規化処理部109は、当該評価対象者フォルマント周波数をパラメータとして、声道長正規化処理を行う。そして、音声処理部110は、所定の母音の音声を受け付けた後に受け付けた音声を処理し、当該音声の評価を行っても良い。

以下、本実施の形態における音声処理装置の具体的な動作について説明する。本具体例において、音声処理装置が語学学習に利用される場合について説明する。

[0059]

[0060]

そして、図示しない手段により、学習した L 種類の音韻 H M M から、学習対象の単語や文章などの音声を構成する 1 以上の音素に対応する H M M を取得し、当該取得した 1 以上の H M M を、音素の順序で連結した教師データを構成する。そして、当該教師データを教師データ格納部 1 0 2 に保持しておく。ここでは、例えば、比較される対象の音声は、単語「right」の音声である。また、ここでは、教師データを発生した者(教師)は、大人である、とする。

[0061]

次に、学習者(評価対象者)が、語学学習の開始の指示である動作開始指示を入力する。かかる指示は、例えば、マウスで所定のボタンを押下することによりなされる。なお、ここでは、学習者は、例えば、子供(5歳から11歳)である、とする。

[0062]

まず、学習者は、母音「う」を発音する、とする。かかる場合、本音声処理装置は、学習に、「う」を発声するように促すことは好適である。「う」を発声するように促すために、音声処理装置は、例えば、「"う"と発声してください。」と画面出力しても良いし、

10

20

30

40

10

20

30

40

50

「"う"と発声してください。」と音声出力しても良い。また、母音「う」は、学習者の評価対象者フォルマント周波数を取得するために好適である。また、本音声処理装置は、第一サンプリング周波数として、「22.05KHz」を保持している、とする。

そして、次に、サンプリング部106は、音声受付部103が受け付けた音声「う」を サンプリングし、「う」の第一音声データを得る。

[0063]

次に、評価対象者フォルマント周波数取得部107は、サンプリング部106が音声「う」をサンプリングして得た第一音声データから、第二フォルマント周波数を取得する。そして、この第二フォルマント周波数を評価対象者フォルマント周波数(Fiとする。今、このFiが「1725Hz」であった、とする。そして、評価対象者フォルマント周波数取得部107は、Fi(1725Hz)を、評価対象者フォルマント周波数格納部108に一時格納する。

[0064]

次に、声道長正規化処理部109は、教師データフォルマント周波数格納部104の教師データフォルマント周波数を読み出す。教師データフォルマント周波数格納部104に格納されている教師データフォルマント周波数は、大人の第二フォルマント周波数であり、今、「1184Hz」である、とする。また、教師データフォルマント周波数は、例えば、教師データを構築する場合に、教師に、例えば、「う」と発声してもらい、当該音声「う」をサンプリング処理した後、取得した第二フォルマント周波数である。

[0065]

なお、図5に、年齢層別、性別ごとの、「う」の第一フォルマント周波数(F1)、第 ニフォルマント周波数(F2)の計測結果を示す。図5により、年齢、性別により、第一 フォルマント周波数(F1)、第二フォルマント周波数(F2)の値が大きく異なること が分る。

[0066]

そして、次に、声道長正規化処理部109は、評価対象者フォルマント周波数「1725 H z 」と教師データフォルマント周波数「1184 H z 」から演算「教師データフォルマント周波数 / 評価対象者フォルマント周波数」を行い、周波数スケール(r)を得る。具体的には、声道長正規化処理部109は、「1184/1725」により、周波数スケール「0.686」を得る。

[0067]

次に、声道長正規化処理部109は、第一サンプリング周波数(Fs)と「r」に基づいて、演算「Fs/r」を行い、第二サンプリング周波数(Fs/r)を得る。ここで、得た第二サンプリング周波数は、「22.05/0.686」により、「32.1」である。そして、声道長正規化処理部109は、第二サンプリング周波数「32.1KHz」を一時格納する。

[0068]

次に、学習者は、学習対象の音声「right」を発音する。そして、音声受付部103は、学習者が発音した音声の入力を受け付ける。なお、音声処理装置は、学習者に「"right"を発音してください。」などを表示、または音声出力するなどして、学習者に「right」の発声を促すことは好適である。

[0069]

次に、サンプリング部106は、受け付けた音声「right」をサンプリング周波数「22.05KHz」でサンプリング処理する。そして、サンプリング部106は、音声「right」の第一音声データを得る。

次に、声道長正規化処理部109は、「right」の第一音声データを第二サンプリング周波数「32.1 KHz」でリサンプリング処理する。そして、声道長正規化処理部109は、第二音声データを得る。

次に、音声処理部110は、第二音声データを、以下のように処理する。

まず、フレーム区分手段1101は、第二音声データを、短時間フレームに区分する。

なお、フレームの間隔は、予め決められている、とする。

[0070]

そして、フレーム音声データ取得手段1102は、フレーム区分手段1101が区分した音声データを、スペクトル分析し、特徴ベクトル系列「 $O=o_1$, o_2 , $\cdot\cdot\cdot$, o_T 」を算出する。なお、Tは、系列長である。ここで、特徴ベクトル系列は、各フレームの特徴ベクトルの集合である。また、特徴ベクトルは、例えば、三角型フィルタを用いたチャネル数24のフィルタバンク出力を離散コサイン変換したMFCCであり、その静的パラメータ、デルタパラメータおよびデルタデルタパワー(39次元)を有する。ま正規化されたパワーとデルタパワーおよびデルタデルタパワー(39次元)を有する。また、スペクトル分析において、ケプストラム平均除去を施すことは好適である。なお、音声分析条件の例を図6の表に示す。なお、音声分析条件は、他の実施の形態における具体例の説明においても同様である。ただし、音声分析条件が、他の条件でも良いことは言うまでもない。また、音声分析の際のサンプリング周波数は、第一サンプリング周波数「205 KHz」である。

[0071]

[0072]

次に、最適状態確率値取得手段 $1\ 1\ 0\ 3\ 2$ は、以下の数式 $1\$ により、最適状態における最適状態確率値($_{t}$ (q_{t} *))を算出する。なお、 $_{t}$ (q_{t} *) は、状態 j の事後確率関数 $_{t}$ (j) の j に q_{t} * を代入した値である。そして、状態 j の事後確率関数 $_{t}$ (j) は、数式 2 を用いて算出される。この確率値($_{t}$ (j))は、 t 番目の特徴ベクトル $_{t}$ が状態 j から生成された事後確率であり、動的計画法を用いて算出される。なお、j は、状態を識別する状態識別子である。

【数1】

$\mathsf{DAP}(\mathsf{t}) = \gamma_{\scriptscriptstyle{+}}(\mathsf{q}_{\scriptscriptstyle{+}}^*)$

数式 1 において、 q_t は、 o_t に対する状態識別子を表す。この確率値(t_t (t_t (t_t))は、 t_t H M M の最尤推定における B a u m - W e 1 c h アルゴリズムの中で表れる占有度数に対応する。

10

20

20

30

40

50

【数2】

$$\gamma_{t}(j) = \frac{\Pr(q_{t}=j, \mathbf{O} \mid \lambda^{\text{all}})}{\Pr(\mathbf{O} \mid \lambda^{\text{all}})}$$

$$= \frac{\Pr(q_{t}=j, \mathbf{O} \mid \lambda^{\text{all}})}{\sum_{k=1}^{N} \Pr(q_{t}=k, \mathbf{O} \mid \lambda^{\text{all}})} \qquad (1)$$

$$= \frac{\alpha_{t}(j) \beta_{t}(j)}{\sum_{k=1}^{N} \alpha_{t}(k) \beta_{t}(k)} \qquad (2)$$

$$\alpha_{t}(j) = Pr(q_{t}=j, \{\mathbf{o}_{1}, \mathbf{o}_{2}, \cdots, \mathbf{o}_{t}\} \mid \lambda^{all})$$

$$\beta_{t}(j) = \Pr(\{\mathbf{o}_{t+1}, \ \mathbf{o}_{t+2}, \ \cdots, \ \mathbf{o}_{T}\} \mid q_{t} = j, \ \lambda^{\text{all}})$$

[0073]

数式 2 において、「 $_{t}$ ($_{j}$)」「 $_{t}$ ($_{j}$)」は、全部の $_{t}$ H M M を用いて、 $_{t}$ forward - $_{t}$ backward アルゴリズムにより算出される。「 $_{t}$ ($_{j}$)」は前向き尤度、「 $_{t}$ ($_{j}$)」は後向き尤度である。 $_{t}$ B a u m - Welch アルゴリズム、 $_{t}$ forward - $_{t}$ backward アルゴリズムは、公知のアルゴリズムであるので、詳細な説明は省略する。

また、数式2において、Nは、全HMMに渡る状態の総数を示す。

[0074]

なお、評定手段1103は、まず最適状態を求め、次に、最適状態の確率値(なお、確率値は、0以上、1以下である。)を求めても良いし、評定手段1103は、まず、全状態の確率値を求め、その後、特徴ベクトル系列の各特徴ベクトルに対する最適状態を求め、当該最適状態に対応する確率値を求めても良い。

[0075]

次に、評定値算出手段11033は、例えば、上記の取得した最適状態確率値と、当該最適状態確率値に対応するフレームの全状態における確率値の総和をパラメータとして表声ができる。かかる場合、もし学習者のセフレーム目に対応する発声が、教式2の(2)教の分子の値が、他の全ての可能な音韻の全ての状態と比較して大きくなり、結果的に近くなの確率値(評定値)が大きくなる。逆にその区間が、教師データが示す発音に近くなる。逆にその区間が、教師データが示す発音に近くなる。がは全ての区間が、教師データが示す発音に近くないような場合は、でははほぼ1/Nに等しくなる。Nは全ての音韻HMMにおける全ての状態の数であるがでははほぼ1/Nに等しくなる。Nは全ての音韻HMMにおける全ての状態の数である。状態における確率値と全ての可能な状態における確率値との比率で定義されている。状態における確率値とのにより多少のスペクトルの変動があったとしても、学習をがって、収音環境等の違いにより多少のスペクトルの変動があったとしても、学習で、いのでで、収音環境等の違いにより多少のスペクトルの変動があったとしても、学習で、では算出手段11033は、最適状態確率値取得手段11032が取得した確率値といい発音をしていれば、その変動が相殺を値取得手段11032が取得した確率値を算出することは、極めて好適である。

[0076]

かかる評定値算出手段11033が算出した評定値(「DAPスコア」とも言う。)を、図7、図8に示す。図7、図8において、横軸は分析フレーム番号、縦軸はスコアを%で表わしたものである。太い破線は音素境界,細い点線は状態境界(いずれもViter

biアルゴリズムで求まったもの)を表わしており、図の上部に音素名を表記している。図7は、アメリカ人男性による英語「right」の発音のDAPスコアを示す。なお、評定値を示すグラフの横軸、縦軸は、後述するグラフにおいても同様である。

[0077]

図8は、日本人男性による英語「right」の発音のDAPスコアを示す。アメリカ人の発音は、日本人の発音と比較して、基本的にスコアが高い。なお、図7において、状態の境界において所々スコアが落ち込んでいることがわかる。

[0078]

そして、出力手段1104は、評定手段1103の評定結果を出力する。具体的には、例えば、出力手段1104は、図9に示すような態様で、評定結果を出力する。つまり、出力手段1104は、各フレームにおける発音の良さを表すスコア(スコアグラフ)として、各フレームの評定値を表示する。その他、出力手段1104は、学習対象の単語の表示(単語表示)、音素要素の表示(音素表示)、教師データの波形の表示(教師波形)、学習者の入力した発音の波形の表示(ユーザ波形)を表示しても良い。なお、図9において、「録音」ボタンを押下すれば、動作開始指示が入力されることとなり、「停止」ボタンを押下すれば、終了指示が入力されることとなる。また、音素要素の表示や波形の表示をする技術は公知技術であるので、その詳細説明を省略する。また、本音声処理装置は、学習対象の単語(図9の「word1」など)や、音素(図9の「p1」など)や、教師波形を出力されるためのデータを予め格納している、とする。

[0079]

また、図9において、フレーム単位以外に、音素単位、単語単位、発声全体の評定結果を表示しても良い。上記の処理において、フレーム単位の評定値を算出するので、単語単位、発声全体の評定結果を得るためには、フレーム単位の1以上の評定値をパラメータとして、単語単位、発声全体の評定値を算出する必要がある。かかる算出式は問わないが、例えば、単語を構成するフレーム単位の1以上の評定値の平均値を単語単位の評定値とする、ことが考えられる。

[0800]

なお、図9において、音声処理装置は、波形表示(教師波形またはユーザ波形)の箇所においてクリックを受け付けると、再生メニューを表示し、音素区間内ではその音素またはその区間が属する単語、波形全体を再生し、単語区間外(無音部)では波形全体のみを再生するようにしても良い。

また、出力手段1104の表示は、図10に示すような態様でも良い。図10において、音素ごとのスコア、単語のスコア、総合スコアが、数字で表示されている。

なお、出力手段1104の表示は、図7、図8のような表示でも良いことは言うまでもない。

[0081]

以上、本実施の形態によれば、ユーザが入力した発音を、教師データに対して、如何に似ているかを示す類似度(評定値)を算出し、出力できる。また、かかる場合、本実施の 形態によれば、個人差、特に声道長の違いに影響を受けない、精度の高い評定ができる。

[0082]

また、本実施の形態によれば、連結された H M M である連結 H M M を用いて最適状態を求め、評定値を算出するので、高速に評定値を求めることができる。したがって、上記の具体例で述べたように、リアルタイムに、フレームごと、音素ごと、単語ごとの評定値を出力できる。また、本実施の形態によれば、動的計画法に基づいた事後確率を確率値として算出するので、さらに高速に評定値を求めることができる。また、本実施の形態によれば、フレームごとに確率値を算出するので、上述したように、フレーム単位だけではなく、または / および音素・単語単位、または / および発声全体の評定結果を出力でき、出力態様の自由度が高い。

[0083]

また、本実施の形態によれば、音声処理装置は、語学学習に利用することを主として説

10

20

30

40

10

20

30

40

50

明したが、物真似練習や、カラオケ評定や、歌唱評定などに利用できる。つまり、本音声処理装置は、比較される対象の音声に関するデータとの類似度を精度良く、高速に評定し、出力でき、そのアプリケーションは問わない。つまり、例えば、本音声処理装置は、カラオケ評価装置であって、音声受付部は、評価対象者の歌声の入力を受け付け、音声処理部は、前記歌声を評価する、という構成でも良い。かかることは、他の実施の形態においても同様である。

[0084]

また、本実施の形態において、音声の入力を受け付けた後または停止ボタン操作後に、スコアリング処理を実行するかどうかをユーザに問い合わせ、スコアリング処理を行うとの指示を受け付けた場合のみ、図10に示すような音素スコア、単語スコア、総合スコアを出力するようにしても良い。

[0085]

また、本実施の形態において、教師データは、比較される対象の音声に関するデータであり、音韻毎の隠れマルコフモデル(HMM)に基づくデータであるとして、主として説明したが、必ずしもHMMに基づくデータである必要はない。教師データは、単一ガウス分布モデルや、確率モデル(GMM:ガウシャンミクスチャモデル)や統計モデルなど、他のモデルに基づくデータでも良い。かかることは、他の実施の形態においても同様である。

[0086]

また、本実施の形態の具体例において、学習者は、母音「う」を発音し、音声処理装置は、かかる音声から第二サンプリング周波数を得た。しかし、学習者は、例えば、母音「あいえお」等、1以上の母音を発音し、かかる母音の音声から、音声処理装置は、第二サンプリング周波数を得ても良い。つまり、第二サンプリング周波数を得るために、学習者が発音する音は「う」に限られない。

[0087]

また、本実施の形態において、音声処理装置が行う下記の処理を、一のDSP(デジタルシグナルプロセッサ)で行っても良い。つまり、本DSPは、第一サンプリング周波数を格納している第一サンプリング周波数格納部と、前記第一サンプリング周波数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリング部と、前記教師データのフォルマント周波数である教師データフォルマント周波数を格納している教師データフォルマント周波数格納部と、前記音声の話者である評価対象者のフォルマント周波数である評価対象者フォルマント周波数を格納している評価対象者フォルマント周波数格納部と、第二サンプリング周波数「前記第一サンプリング周波数/(教師データフォルマント周波数)」で、前記受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理部を具備するデジタルシグナルプロセッサ、である。かかることは、他の実施の形態でも同様である。

[0088]

さらに、本実施の形態における処理は、ソフトウェアで実現しても良い。そして、このソフトウェアをソフトウェアダウンロード等により配布しても良い。また、このソフトウェアをCD-ROMなどの記録媒体に記録して流布しても良い。なお、このことは、本明細書における他の実施の形態においても該当する。なお、本実施の形態における音声処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータに、第一サンプリング周波数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリングステップと、第二サンプリング周波数「第一サンプリング周波数/(教師データフォルマント周波数/評価対象者フォルマント周波数)」で、前記音声受付ステップで受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理ステップと、前記第二音声データを処理する音声処理ステップを実行させるためのプログラム、である。

[0089]

また、上記プログラムにおいて、音声処理ステップは、前記第二音声データを、フレー

ムに区分するフレーム区分ステップと、前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム音声データ取得ステップと、前記教師データと前記1以上のフレーム音声データに基づいて、前記受け付けた音声の評定を行う評定ステップと、前記評定ステップにおける評定結果を出力する出力ステップを具備する、ことは好適である。

[0090]

さらに、上記プログラムにおいて、前記評定ステップは、前記1以上のフレーム音声データのうちの少なくとも一のフレーム音声データに対する最適状態を決定する最適状態決定ステップと、前記最適状態決定ステップで決定した最適状態における確率値を取得する最適状態確率値取得ステップと、前記最適状態確率値取得ステップで取得した確率値をパラメータとして音声の評定値を算出する評定値算出ステップを具備することは好適である

10

(実施の形態2)

[0091]

本実施の形態における音声処理装置は、実施の形態1の音声処理装置と比較して、評定部における評定アルゴリズムが異なる。本実施の形態において、評定値は、最適状態を含む音韻の中の全状態の確率値を発音区間で評価して、算出される。本実施の形態における音声処理装置が算出する事後確率を、実施の形態1におけるDAPに対してt-p-DAPと呼ぶ。

[0092]

20

図11は、本実施の形態における音声処理装置のブロック図である。本音声処理装置は、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部104、第一サンプリング周波数格納部105、サンプリング部106、評価対象者フォルマント周波数取得部107、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部1110、発声催促部1109を具備する

音声処理部 1 1 1 0 は、フレーム区分手段 1 1 0 1、フレーム音声データ取得手段 1 1 0 2、評定手段 1 1 1 0 3、出力手段 1 1 0 4 を具備する。

評定手段11103は、最適状態決定手段11031、発音区間フレーム音韻確率値取得手段111032、評定値算出手段111033を具備する。

30

発音区間フレーム音韻確率値取得手段111032は、最適状態決定手段11031が 決定した各フレームの最適状態を有する音韻全体の状態における1以上の確率値を、発音 区間毎に取得する。

[0093]

評定値算出手段111033は、発音区間フレーム音韻確率値取得手段111032が取得した1以上の発音区間毎の1以上の確率値をパラメータとして音声の評定値を算出する。評定値算出手段111033は、例えば、最適状態決定手段11031が決定した各フレームの最適状態を有する音韻全体の状態における1以上の確率値の総和を、フレーム毎に得て、当該フレーム毎の確率値の総和に基づいて、発音区間毎の確率値の総和の時間平均値を1以上得て、当該1以上の時間平均値をパラメータとして音声の評定値を算出する。

40

[0094]

発音区間フレーム音韻確率値取得手段111032、および評定値算出手段11103 3 は、通常、MPUやメモリ等から実現され得る。発音区間フレーム音韻確率値取得手段 1 1 1 0 3 2 等の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはRO M等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0095]

発声催促部1109は、入力受付部101が、動作開始指示を受け付けた場合、第二サンプリング周波数を算出するために、評価対象者に発声を促す処理を行ったり、評価対象者の発音評定のために発声を促す処理を行ったりする。評価対象者に発声を促す処理は、

10

20

30

40

50

例えば、「~を発音してください。」とディスプレイに表示したり、「~を発音してください。」とスピーカーから音出力したりする処理である。発声催促部1109は、ディスプレイ344やスピーカー等の出力デバイスを含むと考えても含まないと考えても良い。発声催促部1109は、出力デバイスのドライバーソフトまたは、出力デバイスのドライバーソフトと出力デバイス等で実現され得る。

次に、本音声処理装置の動作について図12から図14のフローチャートを用いて説明する。図12等のフローチャートにおいて、図2、図3のフローチャートと異なるステップについてのみ説明する。

(ステップS1201)発声催促部1109は、第二サンプリング周波数算出用の発声を促すために、例えば、母音「う」と発声してください、とディスプレイに表示する。

(ステップS1202)音声受付部103は、音声を受け付けたか否かを判断する。音声を受け付ければステップS1203に行き、音声を受け付けなければステップS213に行く。

[0096]

(ステップS1203)サンプリング部106は、第一サンプリング周波数格納部10 5に格納されている第一サンプリング周波数を読み込み、当該第一サンプリング周波数で 、ステップS1202で受け付けた音声をサンプリングし、第一音声データを得る。

[0097]

(ステップS1204) 声道長正規化処理部109は、ステップS1203で得た第一音声データから、第二サンプリング周波数を得る。かかる第二サンプリング周波数算出処理は、図13のフローチャートを用いて説明する。

(ステップS1205)発声催促部1109は、評定用の発声を促すために、例えば、「right」と発声してください、とディスプレイに表示する。

(ステップS1206)音声受付部103は、音声を受け付けたか否かを判断する。音声を受け付ければステップS1207に行き、音声を受け付けなければステップS213に行く。

[0098]

(ステップS1207)サンプリング部106は、第一サンプリング周波数格納部10 5に格納されている第一サンプリング周波数を読み込み、当該第一サンプリング周波数で 、ステップS1206で受け付けた音声をサンプリングし、第一音声データを得る。

(ステップS1208)声道長正規化処理部109は、ステップS1207で得た第一音声データに対して、ステップS1204で得た第二サンプリング周波数で、リサンプリングし、第二音声データを得る。

(ステップS1209)音声処理部1110は、ステップS1208で得た第二音声データに対して、評定処理を行う。評定処理の詳細は、図14のフローチャートを用いて説明する。ステップS1202に戻る。

なお、図12のフローチャートにおいて、第二サンプリング周波数を算出するための音声と、評定するための音声が同一または包含されていても良い。

[0099]

ステップ S 1 2 0 4 の第二サンプリング周波数算出処理について、図 1 3 のフローチャートを用いて説明する。図 1 3 のフローチャートにおいて、図 3 のフローチャートにおけるステップ S 3 0 1 からステップ S 3 0 5 の処理を行う。

ステップS1209の評定処理について、図14のフローチャートを用いて説明する。 -

[0100]

(ステップS1401)発音区間フレーム音韻確率値取得手段111032は、全フレームの全状態の前向き尤度と後向き尤度を算出する。そして、全フレーム、全状態の確率値を得る。具体的には、発音区間フレーム音韻確率値取得手段111032は、例えば、各特徴ベクトルが対象の状態から生成された事後確率を算出する。この事後確率は、HMMの最尤推定におけるBaum-Welchアルゴリズムの中で現れる占有度数に対応する。Baum-Welchアルゴリズムは、公知のアルゴリズムであるので、説明は省略

する。

(ステップS 1 4 0 2)発音区間フレーム音韻確率値取得手段 1 1 1 0 3 2 は、全フレームの最適状態確率値を算出する。

(ステップS 1 4 0 3) 発音区間フレーム音韻確率値取得手段 1 1 1 0 3 2 は、 j に 1 を代入する。

[0101]

(ステップS1404)発音区間フレーム音韻確率値取得手段111032は、次の評定対象の発音区間である、 j番目の発音区間が存在するか否かを判断する。 j番目の発音区間が存在すればステップS1403に行き、 j番目の発音区間が存在しなければ上位関数にリターンする。

(ステップS 1 4 0 5) 発音区間フレーム音韻確率値取得手段 1 1 1 0 3 2 は、カウンタ k に 1 を代入する。

[0102]

(ステップS1406)発音区間フレーム音韻確率値取得手段111032は、k番目のフレームが、j番目の発音区間に存在するか否かを判断する。k番目のフレームが存在すればステップS1407に行き、k番目のフレームが存在しなければステップS1410に飛ぶ。

(ステップS 1 4 0 7)発音区間フレーム音韻確率値取得手段 1 1 1 0 3 2 は、 k 番目のフレームの最適状態を含む音韻の全ての確率値を取得する。

(ステップS1408)評定値算出手段111033は、ステップS1407で取得した1以上の確率値をパラメータとして、1フレームの音声の評定値を算出する。

(ステップS1409)発音区間フレーム音韻確率値取得手段111032は、 k を 1、インクメントする。ステップS1406に戻る。

[0 1 0 3]

(ステップS1410)評定値算出手段111033は、 j番目の発音区間の評定値を 算出する。評定値算出手段111033は、例えば、最適状態決定手段11031が決定 した各フレームの最適状態を有する音韻全体の状態における1以上の確率値の総和を、フ レーム毎に得て、当該フレーム毎の確率値の総和に基づいて、発音区間の確率値の総和の 時間平均値を、当該発音区間の音声の評定値として算出する。

(ステップS1411)出力手段1104は、ステップS1410で算出した評定値を 出力する。

(ステップS1412)発音区間フレーム音韻確率値取得手段111032は、 j を 1、インクメントする。ステップS1404に戻る。

以下、本実施の形態における音声処理装置の具体的な動作について説明する。本実施の 形態において、評定値の算出アルゴリズムが実施の形態1とは異なるので、その動作を中 心に説明する。

[0104]

まず、学習者(評価対象者)が、語学学習の開始の指示である動作開始指示を入力する。そして、音声処理装置は、当該動作開始指示を受け付け、次に、発声催促部1109は、例えば、「"う"と発声してください。」と画面出力する。

[0105]

なお、ここでも、例えば、学習者は、実施の形態1と同様に子供である。また、教師データを作成するために発声した教師は、ネイティブの大人である、とする。かかることは、他の実施の形態の具体例の記載においても同様である、とする。

そして、評価対象者は、"う"と発声し、音声処理装置は、当該発声から、第二ンプリング周波数「32.1 K H z 」を得る。かかる処理は、実施の形態1において説明した処理と同様である。

[0106]

次に、発声催促部1109は、例えば、「"right"と発声してください。」と画面出力する。そして、学習者は、学習対象の音声「right」を発音する。そして、音声

10

20

30

40

10

20

30

50

受付部103は、学習者が発音した音声の入力を受け付ける。

次に、サンプリング部106は、受け付けた音声「right」をサンプリング周波数「22.05KHz」でサンプリング処理する。そして、サンプリング部106は、「right」の第一音声データを得る。

[0107]

次に、声道長正規化処理部109は、「right」の第一音声データを第二サンプリング周波数「32.1 KHz」でリサンプリング処理する。そして、声道長正規化処理部109は、第二音声データを得る。次に、音声処理部1110は、第二音声データを、以下のように処理する。

まず、フレーム区分手段1101は、「right」の第二音声データを、短時間フレームに区分する。

そして、フレーム音声データ取得手段 1 1 0 2 は、フレーム区分手段 1 1 0 1 が区分した音声データを、スペクトル分析し、特徴ベクトル系列「 $O = O_1$, O_2 , ・・・, O_T 」を算出する。

次に、発音区間フレーム音韻確率値取得手段111032は、各フレームの各状態の事後確率(確率値)を算出する。確率値の算出は、上述した数式1、数式2により算出できる。

[0108]

次に、最適状態決定手段11031は、取得した特徴ベクトル系列を構成する各特徴ベクトルο_tに基づいて、各フレームの最適状態(特徴ベクトルο_tに対する最適状態)を決定する。つまり、最適状態決定手段11031は、最適状態系列を得る。なお、各フレームの各状態の事後確率(確率値)を算出する処理と、最適状態を決定する処理の処理順序は問わない。

[0109]

次に、発音区間フレーム音韻確率値取得手段111032は、発音区間ごとに、当該発音区間に含まれる各フレームの最適状態を含む音韻の全ての確率値を取得する。そして、評定値算出手段111033は、フレーム毎に算出する。そして、評定値算出手段111033は、フレーム毎に算出する。して、評定値算出手段111033は、フレーム毎に算出された確率値の総和を、発音区間毎に時間平均し、発音区間毎の評定値を算出する。具体的には、評定値算出手段111033は、数式3により評定値を算出する。数式3において、p-DAP()は、各フレームにおける、すべての音韻の中で最適な音韻の事後確率(確率値)を表すように算出される評定値であり、数式4で算出され得る。なお、数式3のt-p-DAPは、最適状態を含む音韻の中の全状態の確率値を発音区間で評価して、算出される評定値である。また、数式3において、(♀t゜))は、状態♀t゜を含むHMMが有する全状態識別子の集合である。

【数3】

$$t-p-DAP(m) \stackrel{\sum_{\tau} \in T(q_{t}^{*})}{= |T(q_{t}^{*})|}$$

【数4】

$$\text{p-DAP}\left(t\right) \stackrel{\Delta}{=} \sum_{j \in P(q_{t}^{*})} \gamma_{t} \left(j\right)$$

[0110]

かかる評定値算出手段111033が算出した評定値(「t-p-DAPスコア」とも

言う。)を、図15の表に示す。図15において、アメリカ人男性と日本人男性の評定結果を示す。PhonemeおよびWordは,t-p-DAPにおける時間平均の範囲を示す。ここでは、DAPの代わりにp-DAPの時間平均を採用したものである。図15において、アメリカ人男性の発音の評定値が日本人男性の発音の評定値より高く、良好な評定結果が得られている。

そして、出力手段1104は、算出した発音区間ごと(ここでは、音素毎)の評定値を 、順次出力する。かかる出力例は、図16である。

[0111]

以上、本実施の形態によれば、ユーザが入力した発音を、教師データに対して、如何に似ているかを示す類似度(評定値)を算出し、出力できる。また、かかる場合、本実施の形態によれば、個人差、特に声道長の違いに影響を受けない、精度の高い評定ができる。

[0112]

また、本実施の形態によれば、連結された H M M である連結 H M M を用いて最適状態を求め、評定値を算出するので、高速に評定値を求めることができる。したがって、上記の具体例で述べたように、リアルタイムに、発音区間ごとの評定値を出力できる。また、本実施の形態によれば、動的計画法に基づいた事後確率を確率値として算出するので、さらに高速に評定値を求めることができる。

[0113]

また、本実施の形態によれば、評定値を、発音区間の単位で算出でき、実施の形態 1 におけるような状態単位の DAPと比較して、本来、測定したい類似度(発音区間の類似度)を精度良く、安定して求めることができる。

[0114]

さらに、本実施の形態における音声処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータに、第一サンプリング周波数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリングステップと、第二サンプリング周波数「第一サンプリング周波数 / (教師データフォルマント周波数 / 評価対象者フォルマント周波数)」で、前記音声受付ステップで受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理ステップと、前記第二音声データを処理する音声処理ステップを実行させるためのプログラム、である

[0115]

また、上記プログラムにおいて、音声処理ステップは、前記第二音声データを、フレームに区分するフレーム区分ステップと、前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム音声データ取得ステップと、前記教師データと前記1以上のフレーム音声データに基づいて、前記受け付けた音声の評定を行う評定ステップと、前記評定ステップにおける評定結果を出力する出力ステップを具備する、ことは好適である。

[0116]

また、上記プログラムにおいて、前記評定ステップは、前記1以上のフレーム音声データの最適状態を決定する最適状態決定ステップと、前記最適状態決定ステップで決定した各フレームの最適状態を有する音韻全体の状態における1以上の確率値を、発音区間毎に取得する発音区間フレーム音韻確率値取得ステップと、前記発音区間フレーム音韻確率値取得ステップで取得した1以上の発音区間毎の1以上の確率値をパラメータとして音声の評定値を算出する評定値算出ステップを具備する、ことは好適である。

(実施の形態3)

[0117]

本実施の形態において、比較対象の音声と入力音声の類似度を精度高く評定できる音声処理装置について説明する。特に、本音声処理装置は、無音区間を検知し、無音区間を考慮した類似度評定が可能な音声処理装置である。

[0118]

40

30

20

10

10

20

30

40

50

図17は、本実施の形態における音声処理装置のブロック図である。本音声処理装置は、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部104、第一サンプリング周波数格納部105、サンプリング部106、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部1710、発声催促部1109を具備する

音声処理部 1 7 1 0 は、フレーム区分手段 1 1 0 1、フレーム音声データ取得手段 1 1 0 2、特殊音声検知手段 1 7 1 0 1、評定手段 1 7 1 0 3、出力手段 1 1 0 4 を具備する

特殊音声検知手段 1 7 1 0 1 は、無音データ格納手段 1 7 1 0 1 1、無音区間検出手段 1 7 1 0 1 2 を具備する。

評定手段 1 7 1 0 3 は、最適状態決定手段 1 1 0 3 1、最適状態確率値取得手段 1 1 0 3 2、評定値算出手段 1 7 1 0 3 3 を具備する。

[0119]

特殊音声検知手段17101は、フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する。なお、ここで特殊な音声は、無音も含む。また、特殊音声検知手段17101は、例えば、フレームの最適状態の確率値を、ある音素区間において取得し、ある音素区間の1以上の確率値の総和が所定の値より低い場合(想定されている音素ではない、と判断できる場合)、当該音素区間において特殊な音声が入力されたと、検知する。かかる検知の具体的なアルゴリズムの例は後述する。特殊音声検知手段17101の処理手順は、通常、MPUやメモリ等から実現され得る。特殊音声検知手段17101の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0120]

無音データ格納手段 1 7 1 0 1 1 は、無音を示すデータであり、 H M M に基づくデータである無音データを格納している。無音データ格納手段 1 7 1 0 1 1 は、不揮発性の記録 媒体が好適であるが、揮発性の記録媒体でも実現可能である。

[0121]

無音区間検出手段171012は、フレーム音声データ取得手段1102が取得したフレーム音声データ、および無音データ格納手段171011の無音データに基づいて、無音の区間を検出する。無音区間検出手段171012は、フレーム音声データ取得手段1102が取得したフレーム音声データと無音データの類似度が所定の値以上である場合に、当該フレーム音声データは無音区間のデータであると判断しても良い。また、無音区間検出手段171012は、下記で述べる最適状態確率値取得手段11032が取得したフレーム音声データと無音データの類似度が所定の値以上である場合に、当該フレーム音声データは無音区間のデータであると判断しても良い。無音区間検出手段171012は、通常、MPUやメモリ等から実現され得る。無音区間検出手段171012の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0122]

評定手段17103は、教師データと入力音声データと特殊音声検知手段17101における検知結果に基づいて、音声受付部103が受け付けた音声の評定を行う。「特殊音声検知手段17101における検知結果に基づく」とは、例えば、特殊音声検知手段17101が無音を検知した場合、当該無音の区間を無視することである。また、「特殊音声検知手段17101における検知結果に基づく」とは、例えば、特殊音声検知手段17101が置換や脱落などを検知した場合、当該置換や脱落などの検知により、評定値を所定数値分、減じて、評定値を算出することである。評定手段17103は、通常、MPUやメモリ等から実現され得る。評定手段17103の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア

(専用回路)で実現しても良い。

[0123]

評定値算出手段171033は、無音区間検出手段171012が検出した無音区間を除いて、かつ最適状態確率値取得手段11032が取得した確率値をパラメータとして音声の評定値を算出する。なお、評定値算出手段171033は、上記確率値を如何に利用して、評定値を算出するかは問わない。評定値算出手段171033は、例えば、最適状態確率値取得手段11032が取得した確率値と、当該確率値に対応するフレームの全状態における確率値の総和をパラメータとして音声の評定値を算出する。評定値算出手段21023は、ここでは、通常、無音区間検出手段171012が検出した無音区間を除いて、フレームごとに評定値を算出する。なお、評定値算出手段171033は、かならずしも無音区間を除いて、評定値を算出する必要はない。評定値算出手段171033は、無音区間の影響を少なくするように評定値を算出しても良い。評定値算出手段171033は、通常、MPUやメモリ等から実現され得る。評定値算出手段171033の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0124]

次に、音声処理装置の動作について図 1 8、図 1 9 のフローチャートを用いて説明する。なお、図 1 8 のフローチャートは、図 1 2 のフローチャートと比較して、ステップS 1 8 0 1 の評定処理のみが異なるので、図 1 8 のフローチャートの説明は省略する。ステップS 1 8 0 1 の評定処理の詳細について、図 1 9 のフローチャートを用いて説明する。

[0 1 2 5]

(ステップS1901)評定手段17103は、DAPの評定値を算出する。DAPの評定値を算出するアルゴリズムは、実施の形態1で説明済みであるので、詳細な説明は省略する。DAPの評定値を算出する処理は、図2のフローチャートの、ステップS211からS214の処理により行う。

[0126]

(ステップS1902)特殊音声検知手段17101は、ステップS1901で算出した値が、所定の値より低いか否かを判断する。所定の値より低ければステップS1903に行き、所定の値より低くなければステップS1906に飛ぶ。

(ステップS 1 9 0 3) 無音区間検出手段 1 7 1 0 1 2 は、無音データと全教師データの確率値を取得する。

[0127]

(ステップS1904)無音区間検出手段171012は、ステップS1903で取得した確率値の中で、無音データの確率値が最も高いか否かを判断する。無音データの確率値が最も高ければ(かかる場合、無音の区間であると判断する)ステップS1905に行き、無音データの確率値が最も高くなければステップS1906に行く。

(ステップS1905)無音区間検出手段171012は、カウンタiを1、インクリメントする。ステップS208に戻る。

(ステップS1906)出力手段1104は、ステップS1901で算出した評定値を 出力する。

[0128]

なお、図19のフローチャートにおいて、出力手段1104は、無音区間と判定した区間の評定値は出力しなかった(無音区間を無視した)が、特殊音声が検知された区間が無音区間である旨を明示したり、無音区間が存在する旨を明示したりする態様で出力しても良い。また、評定値算出手段171033は、発音区間や、それ以上の単位のスコアを算出する場合に、無音区間の評定値を無視して、スコアを算出することが好適であるが、無音区間の評定値の影響を、例えば、1/10にして、発音区間や発音全体のスコアを算出するなどしても良い。評定手段17103は、教師データと入力音声データと特殊音声検知手段17101における検知結果に基づいて、音声受付部103が受け付けた音声の評定を行えばよい。

10

20

30

[0129]

また、図19のフローチャートにおいて、特殊音声検知手段17101は、i番目のフレーム音声データのDAPスコアに基づいて特殊音声を検知したが、例えば、実施の形態2で説明したt-p-DAPスコアに基づいて特殊音声を検知しても良い。

[0130]

以下、本実施の形態における音声処理装置の具体的な動作について説明する。本実施の 形態において、無音区間を考慮して評定値を算出するので、評定値の算出アルゴリズムが 実施の形態 1、実施の形態 2 とは異なる。そこで、その異なる処理を中心に説明する。

まず、学習者(評価対象者)が、語学学習の開始の指示である動作開始指示を入力する。そして、音声処理装置は、当該動作開始指示を受け付け、次に、例えば、「"う"と発声してください。」と画面出力する。

そして、評価対象者は、"う"と発声し、音声処理装置は、当該発声から、第二ンプリング周波数「32.1 K H z 」を得る。かかる処理は、実施の形態1等において説明した処理と同様である。

[0131]

次に、発声催促部1109は、例えば、「"right"と発声してください。」と画面出力する。そして、学習者は、学習対象の音声「right」を発音する。そして、音声受付部103は、学習者が発音した音声の入力を受け付ける。

次に、サンプリング部106は、受け付けた音声「right」をサンプリング周波数「22.05 K H z」でサンプリング処理する。そして、サンプリング部106は、「right」の第一音声データを得る。

[0132]

次に、声道長正規化処理部109は、「right」の第一音声データを、第二サンプリング周波数「32.1 K H z」でリサンプリング処理する。そして、声道長正規化処理部109は、第二音声データを得る。次に、音声処理部<u>1710</u>は、第二音声データを、以下のように処理する。

まず、フレーム区分手段1101は、「right」の第二音声データを、短時間フレームに区分する。

そして、フレーム音声データ取得手段 1 1 0 2 は、フレーム区分手段 1 1 0 1 が区分した音声データを、スペクトル分析 し、特徴ベクトル系列「 $O=o_1$, o_2 , ・・・, o_T 」を算出する。

次に、最適状態決定手段 1 1 0 3 1 は、取得した特徴ベクトル系列を構成する各特徴ベクトル $_{\rm t}$ に基づいて、所定のフレームの最適状態(特徴ベクトル $_{\rm t}$ に対する最適状態)を決定する。

次に、最適状態確率値取得手段11032は、上述した数式1、2により、最適状態における確率値を算出する。

[0133]

次に、評定値算出手段171033は、例えば、最適状態決定手段11031が決定した最適状態を有する音韻全体の状態における1以上の確率値を取得し、当該1以上の確率値の総和をパラメータとして音声の評定値を算出する。つまり、評定値算出手段171033は、例えば、DAPスコアをフレーム毎に算出する。

[0134]

そして、特殊音声検知手段17101は、算出されたフレームに対応する評定値(DAPスコア)を用いて、特殊な音声が入力されたか否かを判断する。具体的には、特殊音声検知手段17101は、例えば、評価対象のフレームに対して算出された評定値が、所定の数値より低ければ、特殊な音声が入力された、と判断する。なお、特殊音声検知手段17101は、一のフレームに対応する評定値が小さいからといって、直ちに特殊な音声が入力された、と判断する必要はない。つまり、特殊音声検知手段17101は、フレームに対応する評定値が小さいフレームが所定の数以上、連続する場合に、当該連続するフレーム群に対応する区間が特殊な音声が入力された区間と判断しても良い。

10

20

30

40

[0135]

特殊音声検知手段17101が、特殊音声を検知する場合について説明する図を図20に示す。図20(a)の縦軸は、DAPスコアであり、横軸はフレームを示す。図20(a)において、(V)は、Viterbiアライメントを示す。図20(a)において、網掛けのフレーム群のおけるDAPスコアは、所定の値より低く、特殊音声の区間である、と判断される。

[0 1 3 6]

次に、特殊な音声が入力された、と判断した場合、無音区間検出手段171012は、無音データ格納手段171011から無音データを取得し、当該フレーム群と無音データとの類似度を算定し、類似度が所定値以上であれば当該フレーム群に対応する音声データが、無音データであると判断する。図20(b)は、無音データとの比較の結果、当該無音データとの類似度を示す事後確率の値(「DAPスコア」)が高いことを示す。その結果、無音区間検出手段171012は、当該特殊音声の区間は、無音区間である、と判断する。なお、図20(a)において、網掛けのフレーム群のおけるDAPスコアは、所定の値より低く、特殊音声の区間である、と判断され、かつ、無音データとの比較の結果、DAPスコアが低い場合には、無音区間ではない、と判断される。そして、かかる区間において、例えば、単に、発音が上手くなく、低い評定値が出力される。なお、図20(a)に示しているように、通常、無音区間は、第一のワード(「word1」)の最終音素の後半部、および第一のワードに続く第二のワード(「word2」)の第一音素の前半部のスコアが低い。

そして、出力手段1104は、出力する評定値から、無音データの区間の評定値を考慮 しないように、無視する。

そして、出力手段1104は、各フレームに対応する評定値を出力する。この場合、例 えば、無音データの区間の評定値は、出力されない。

かかる評定値の出力態様例は、例えば、図9、図10である。

なお、出力手段1104が行う出力は、無音区間の存在を示すだけの出力でも良い。

[0137]

以上、本実施の形態によれば、ユーザが入力した発音を、教師データに対して、如何に似ているかを示す類似度(評定値)を算出し、出力できる。また、かかる場合、本実施の形態によれば、個人差、特に声道長の違いに影響を受けない、精度の高い評定ができる。さらに、本音声処理装置は、無音区間を考慮して類似度を評定するので、極めて正確な評定結果が得られる。

[0138]

なお、無音区間のデータは、無視して評定結果を算出することは好適である。ただし、本実施の形態において、例えば、無音区間の評価の影響を他の区間と比較して少なくするなど、無視する以外の方法で、無音区間のデータを考慮して、評定値を出力しても良い。 【0139】

また、本実施の形態の具体例によれば、DAPスコアを用いて、評定値を算出したが、無音の区間を考慮して評定値を算出すれば良く、上述した他のアルゴリズム(t - p - D AP等)、または、本明細書では述べていない他のアルゴリズムにより評定値を算出しても良い。つまり、本実施の形態によれば、教師データと入力音声データと特殊音声検知手段における検知結果に基づいて、音声受付部が受け付けた音声の評定を行い、特に、無音データを考慮して、評定値を算出すれば良い。

また、本実施の形態によれば、まず、DAPスコアが低い区間を検出してから、無音区間の検出をした。しかし、DAPスコアが低い区間を検出せずに、無音データとの比較により、無音区間を検出しても良い。

[0140]

さらに、本実施の形態における音声処理装置を実現するソフトウェアは、以下のような プログラムである。つまり、このプログラムは、コンピュータに、第一サンプリング周波 数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリングステッ

10

20

30

40

プと、第二サンプリング周波数「第一サンプリング周波数/(教師データフォルマント周波数/評価対象者フォルマント周波数)」で、前記音声受付ステップで受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理ステップと、前記第二音声データを処理する音声処理ステップを実行させるためのプログラム、である

[0141]

また、上記プログラムにおいて、音声処理ステップは、前記第二音声データを、フレームに区分するフレーム区分ステップと、前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム音声データ取得ステップと、前記フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する特殊音声検知ステップと、教師データと前記入力音声データと前記特殊音声検知ステップにおける検知結果に基づいて、前記受け付けた音声の評定を行う評定ステップと、前記評定ステップにおける評定結果を出力する出力ステップを具備する、ことは好適である。

また、上記プログラムにおいて、特殊音声検知ステップは、無音を示すHMMに基づく データである無音データと、前記入力音声データに基づいて、無音の区間を検出する、こ とは好適である。

(実施の形態4)

[0142]

本実施の形態において、入力音声において、特殊音声を検知し、比較対象の音声と入力音声の類似度を精度高く評定できる音声処理装置について説明する。特に、本音声処理装置は、音韻の挿入を検知できる音声処理装置である。

[0 1 4 3]

図21は、本実施の形態における音声処理装置のブロック図である。本音声処理装置は、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部104、第一サンプリング周波数格納部105、サンプリング部106、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部2110、発声催促部1109を具備する

[0144]

音声処理部2110は、フレーム区分手段1101、フレーム音声データ取得手段1102、特殊音声検知手段21101、評定手段21103、出力手段21104を具備する。なお、評定手段21103は、最適状態決定手段11031、最適状態確率値取得手段11032を具備する。

[0145]

特殊音声検知手段 2 1 1 0 1 は、一の音素の後半部および当該音素の次の音素の前半部の評定値が所定の条件を満たすことを検知する。後半部、および前半部の長さは問わない。特殊音声検知手段 2 1 1 0 1 は、通常、MPUやメモリ等から実現され得る。特殊音声検知手段 2 1 1 0 1 の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0146]

評定手段21103は、特殊音声検知手段21101が所定の条件を満たすことを検知した場合に、少なくとも音素の挿入があった旨を示す評定結果を構成する。なお、評定手段21103は、実施の形態3で述べたアルゴリズムにより、特殊音声検知手段21101が所定の条件を満たすことを検知した区間に無音が挿入されたか否かを判断し、無音が挿入されていない場合に、他の音素が挿入されたと検知しても良い。また、評定手段21103は、無音が挿入されていない場合に、他の音韻HMMに対する確率値を算出し、所定の値より高い確率値を得た音韻が挿入された、との評定結果を得ても良い。なお、実施の形態3で述べた無音区間の検知は、無音音素の挿入の検知である、とも言える。評定手段21103の処理

10

20

30

40

手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0147]

出力手段21104は、評定手段21103における評定結果を出力する。ここでの評定結果は、音素の挿入があった旨を示す評定結果を含む。また、評定結果は、音素の挿入があった場合に、所定数値分、減じられて算出された評定値(スコア)のみでも良い。また、評定結果は、音素の挿入があった旨、および評定値(スコア)の両方であっても良い。なお、教師データにおいて想定されていない音素の挿入を検知した場合、通常、評定値は低くなる。ここで、出力とは、ディスプレイへの表示、プリンタへの印字、音出力、外部の装置への送信、記録媒体への蓄積等を含む概念である。出力手段21104は、ディスプレイやスピーカー等の出力デバイスを含むと考えても含まないと考えても良い。出力手段21104は、出力デバイスのドライバーソフトまたは、出力デバイスのドライバーソフトと出力デバイス等で実現され得る。

[0148]

次に、音声処理装置の動作について、図22、図23のフローチャートを用いて説明する。なお、図22のフローチャートは、図12のフローチャートと比較して、ステップS2201の評定処理のみが異なるので、図22のフローチャートの説明は省略する。ステップS2201の評定処理の詳細について、図23のフローチャートを用いて説明する。図23のフローチャートにおいて、図2、図19のフローチャートの処理と同様の処理については、その説明を省略する。

[0149]

(ステップS2301)特殊音声検知手段21101は、フレームに対応するデータを一時的に蓄積するバッファにデータが格納されているか否かを判断する。なお、格納されているデータは、ステップS1902で、所定の値より低い評定値と評価されたフレーム音声データ、または当該フレーム音声データから取得できるデータである。データが格納されていればステップS2307に行き、データが格納されていなければ上位関数にリターンする。

[0150]

(ステップS2302)特殊音声検知手段21101は、バッファにデータが格納されているか否かを判断する。データが格納されていればステップS2307に行き、データが格納されていなければステップステップS2303に行く。

(ステップS2303)出力手段21104は、ステップS1901で算出した評定値を出力する。

(ステップS2304)特殊音声検知手段21101は、カウンタiを1、インクリメントする。ステップS208に戻る。

(ステップS2305)特殊音声検知手段21101は、バッファに、所定の値より低い評定値と評価されたフレーム音声データ、または当該フレーム音声データから取得できるデータを一時蓄積する。

(ステップ S 2 3 0 6) 特殊音声検知手段 2 1 1 0 1 は、カウンタ i を 1 、インクリメントする。ステップ S 2 0 8 に戻る。

(ステップS2307)特殊音声検知手段21101は、カウンタjに1を代入する。 【0151】

(ステップS2308)特殊音声検知手段21101は、 j 番目のデータが、バッファに存在するか否かを判断する。 j 番目のデータが存在すればステップS2309に行き、 j 番目のデータが存在しなければステップS2315に飛ぶ。

(ステップS2309)特殊音声検知手段21101は、 j番目のデータに対応する最適状態の音素を取得する。

(ステップS2310)特殊音声検知手段21101は、 j番目のデータに対する全教師データの確率値を算出し、最大の確率値を持つ音素を取得する。

[0152]

10

20

30

(ステップS2311)特殊音声検知手段21101は、ステップS2309で取得した音素とステップS2310で取得した音素が異なる音素であるか否かを判断する。異なる音素であればステップS2312に行き、異なる音素でなければステップS2314に飛ぶ。

(ステップS2312)評定手段21103は、音素の挿入があった旨を示す評定結果 を構成する。

(ステップS2313)特殊音声検知手段21101は、カウンタjを1、インクリメントする。ステップS2308に戻る。

(ステップS2314)出力手段21104は、バッファ中の全データに対応する全評定値を出力する。ここで、全評定値とは、例えば、フレーム毎のDAPスコアである。ステップS2313に行く。

[0153]

(ステップS2315)出力手段21104は、評定結果に「挿入の旨」の情報が入っているか否かを判断する。「挿入の旨」の情報が入っていればステップS2316に行き、「挿入の旨」の情報が入っていなければステップS2317に行く。

(ステップS2316)出力手段21104は、評定結果を出力する。

(ステップS2317)出力手段21104は、バッファをクリアする。ステップS2 08に戻る。

[0154]

なお、図23のフローチャートにおいて、評定値の低いフレームが2つの音素に渡って存在すれば、音素の挿入があったと判断した。つまり、一の音素の後半部(少なくとも最終フレーム)および当該音素の次の音素の第一フレームの評定値が所定値より低い場合に、音素の挿入があったと判断した。しかし、図23のフローチャートにおいて、一の音素の所定区間以上の後半部、および当該音素の次の音素の所定区間以上の前半部の評定値が所定値よりすべて低い場合に、音素の挿入があったと判断するようにしても良い。

以下、本実施の形態における音声処理装置の具体的な動作について説明する。本実施の 形態において、音素の挿入の検知を行う処理が実施の形態3等とは異なる。そこで、その 異なる処理を中心に説明する。

まず、学習者(評価対象者)が、語学学習の開始の指示である動作開始指示を入力する。そして、音声処理装置は、当該動作開始指示を受け付け、次に、例えば、「"あ"と発声してください。」と画面出力する。

そして、学習者は、"あ"と発声し、音声処理装置は、当該発声から、第二ンプリング周波数「32.1KHz」を得る。かかる処理は、実施の形態1等において説明した処理と同様である。

[0155]

次に、発声催促部1109は、例えば、「"right"と発声してください。」と画面出力する。そして、学習者は、学習対象の音声「right」を発音する。そして、音声受付部103は、学習者が発音した音声の入力を受け付ける。

次に、サンプリング部106は、受け付けた音声「right」をサンプリング周波数「22.05KHz」でサンプリング処理する。そして、サンプリング部106は、「right」の第一音声データを得る。

[0156]

次に、声道長正規化処理部109は、「right」の第一音声データを第二サンプリング周波数「32.1 KHz」でリサンプリング処理する。そして、声道長正規化処理部109は、第二音声データを得る。次に、音声処理部2110は、第二音声データを、以下のように処理する。

まず、フレーム区分手段1101は、「right」の第二音声データを、短時間フレームに区分する。

そして、フレーム音声データ取得手段 1 1 0 2 は、フレーム区分手段 1 1 0 1 が区分した音声データを、スペクトル分析し、特徴ベクトル系列「 $O = O_1$, O_2 , ・・・, O_T

10

20

30

40

」を算出する。

次に、評定手段21103の最適状態決定手段11031は、取得した特徴ベクトル系列を構成する各特徴ベクトルo_tに基づいて、所定のフレームの最適状態(特徴ベクトルo_tに対する最適状態)を決定する。

次に、最適状態確率値取得手段11032は、上述した数式1、2により、最適状態における確率値を算出する。

[0 1 5 7]

次に、評定手段21103は、例えば、最適状態決定手段11031が決定した最適状態を有する音韻全体の状態における1以上の確率値を取得し、当該1以上の確率値の総和をパラメータとして音声の評定値を算出する。つまり、評定手段21103は、例えば、DAPスコアをフレーム毎に算出する。ここで、算出するスコアは、上述したt-p-DAPスコア等でも良い。

[0158]

そして、特殊音声検知手段 2 1 1 0 1 は、算出されたフレームに対応する評定値を用いて、特殊な音声が入力されたか否かを判断する。つまり、評定値(例えば、DAPスコア)が、所定の値より低い区間が存在するか否かを判断する。

[0159]

次に、特殊音声検知手段21101は、図24に示すように、評定値(例えば、DAP スコア)が、所定の値より低い区間が、2つの音素に跨っているか否かを判断し、2つの 音素に跨がっていれば、当該区間に音素が挿入された、と判断する。なお、かかる場合の 詳細なアルゴリズムの例は、図23で説明した。また、図24において、斜線部が、予期 しない音素が挿入された区間である。

[0160]

次に、評定手段 2 1 1 0 3 は、音素の挿入があった旨を示す評定結果(例えば、「予期しない音素が挿入されました。」)を構成する。なお、予期しない音素が挿入された場合、評定手段 2 1 1 0 3 は、例えば、所定数値分、減じて、評定値を算出することは好適である。そして、出力手段 2 1 1 0 4 は、構成した評定結果(評定値を含んでも良い)を出力する。図 2 5 は、評定結果の出力例である。なお、出力手段 2 1 1 0 4 は、通常の入力音声に対しては、上述したように評定値を出力することが好適である。

[0161]

以上、本実施の形態によれば、ユーザが入力した発音を、教師データに対して、如何に似ているかを示す類似度(評定値)を算出し、出力できる。また、かかる場合、本実施の形態によれば、個人差、特に声道長の違いに影響を受けない、精度の高い評定ができる。さらに、本音声処理装置は、特殊音声、特に、予期せぬ音素の挿入を検知できるので、極めて精度の高い評定結果が得られる。

[0162]

なお、本実施の形態において、音素の挿入を検知できれば良く、評定値の算出アルゴリズムは問わない。評定値の算出アルゴリズムは、上述したアルゴリズム(DAP、t-p-DAP)でも良く、または、本明細書では述べていない他のアルゴリズムでも良い。

[0163]

さらに、本実施の形態における音声処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータに、第一サンプリング周波数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリングステップと、第二サンプリング周波数「第一サンプリング周波数/(教師データフォルマント周波数/評価対象者フォルマント周波数)」で、前記音声受付ステップで受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理ステップと、前記第二音声データを処理する音声処理ステップを実行させるためのプログラム、である

[0164]

また、上記プログラムにおいて、音声処理ステップは、前記第二音声データを、フレー

10

20

30

40

ムに区分するフレーム区分ステップと、前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム音声データ取得ステップと、前記フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する特殊音声検知ステップと、教師データと前記入力音声データと前記特殊音声検知ステップにおける検知結果に基づいて、前記受け付けた音声の評定を行う評定ステップと、前記評定ステップにおける評定結果を出力する出力ステップを具備する、ことは好適である。

また、上記プログラムにおいて、特殊音声検知ステップは、一の音素の後半部および当該音素の次の音素の前半部の評定値が所定の条件を満たすことを検知する、ことは好適である。

(実施の形態5)

[0165]

本実施の形態において、入力音声において、特殊音声を検知し、比較対象の音声と入力音声の類似度を精度高く評定できる音声処理装置について説明する。特に、本音声処理装置は、音韻の置換を検知できる音声処理装置である。

[0166]

図26は、本実施の形態における音声処理装置のブロック図である。本音声処理装置は、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部104、第一サンプリング周波数格納部105、サンプリング部106、評価対象者フォルマント周波数格納部107、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部2610、発声催促部1109を具備する

[0167]

音声処理部2610は、フレーム区分手段1101、フレーム音声データ取得手段1102、特殊音声検知手段26101、評定手段26103、出力手段21104を具備する。なお、評定手段26103は、最適状態決定手段11031、最適状態確率値取得手段11032を具備する。なお、評定手段26103は、最適状態決定手段11031、最適状態確率値取得手段11032を具備する。

[0168]

音声処理部2610は、第二音声データを処理する。音声処理部2610は、フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する特殊音声検知手段26101を具備する。音声処理部2610は、通常、MPUやメモリ等から実現され得る。音声処理部2610の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0169]

特殊音声検知手段 2 6 1 0 1 は、一の音素の評定値が所定の値より低いことを検知する。また、特殊音声検知手段 2 6 1 0 1 は、一の音素の評定値が所定の値より低く、かつ当該音素の直前の音素および当該音素の直後の音素の評定値が所定の値より高いことをも検知しても良い。また、特殊音声検知手段 2 6 1 0 1 は、一の音素の評定値が所定の値より低く、かつ、想定していない音素のHMMに基づいて算出された評定値が所定の値より高いことを検知しても良い。つまり、特殊音声検知手段 2 6 1 0 1 は、所定のアルゴリズムで、音韻の置換を検知できれば良い。そのアルゴリズムは種々考えられる。特殊音声検知手段 2 6 1 0 1 は、通常、MPUやメモリ等から実現され得る。特殊音声検知手段 2 6 1 0 1 の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0170]

評定手段26103は、特殊音声検知手段26101が所定の条件を満たすことを検知した場合に、少なくとも音素の置換があった旨を示す評定結果を構成する。評定手段26103は、音素の置換があった場合に、所定数値分、減じられて算出された評定値(スコア)を算出しても良い。評定手段26103は、通常、MPUやメモリ等から実現され得

10

20

40

30

る。評定手段26103の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0171]

次に、音声処理装置の動作について、図27、図28のフローチャートを用いて説明する。なお、図27のフローチャートは、図12のフローチャートと比較して、ステップS2701の評定処理のみが異なるので、図27のフローチャートの説明は省略する。ステップS2701の評定処理の詳細について、図28のフローチャートを用いて説明する。図28のフローチャートにおいて、図2、図19、図23のフローチャートの処理と同様の処理については、その説明を省略する。

[0172]

(ステップS2801)特殊音声検知手段26101は、バッファに蓄積されているデータに対応するフレーム音声データ群が一の音素に対応するか否かを判断する。一の音素であればステップS2802に行き、一の音素でなければステップS2810に行く。

[0 1 7 3]

(ステップS2802)特殊音声検知手段26101は、バッファに蓄積されているデータに対応するフレーム音声データ群の音素の直前の音素の評定値を算出する。かかる評定値は、例えば、上述したDAPスコアである。なお、直前の音素とは、現在評定中の音素に対して直前の音素である。音素の区切りは、Viterbiアルゴリズムにより算出できる。

[0174]

(ステップS2803)特殊音声検知手段26101は、ステップS2802で算出した評定値が所定の値以上であるか否かを判断する。所定の値以上であればステップS28 04に行き、所定の値より小さければステップS2810に行く。

(ステップS2804)特殊音声検知手段26101は、直後の音素の評定値を算出する。かかる評定値は、例えば、上述したDAPスコアである。直後の音素とは、現在評定中の音素に対して直後の音素である。

[0175]

(ステップS2805)特殊音声検知手段26101は、ステップS2804で算出した評定値が所定の値以上であるか否かを判断する。所定の値以上であればステップS28 06に行き、所定の値より小さければステップS2810に行く。

[0176]

(ステップS2806)特殊音声検知手段26101は、予め格納されている音韻HMM(予期する音韻のHMMは除く)の中で、所定の値以上の評定値が得られる音韻HMMが一つ存在するか否かを判断する。所定の値以上の評定値が得られる音韻HMMが存在すればステップS2807に行き、所定の値以上の評定値が得られる音韻HMMが存在しなければステップS2810に行く。なお、予め格納されている音韻HMMは、通常、すべての音韻に対する多数の音韻HMMである。なお、本ステップにおいて、予め格納されている音韻HMMの確率値を算出し、最大の確率値を持つ音素を取得し、当該音素と最適状態の音素が異なるか否かを判断し、異なる場合に音素の置換があったと判断しても良い。

(ステップS2807)評定手段26103は、音素の置換があった旨を示す評定結果 を構成する。

(ステップS2808)出力手段21104は、ステップS2807で構成した評定結果を出力する。

(ステップS2809)出力手段21104は、バッファをクリアする。ステップS2 08に戻る。

(ステップ S 2 8 1 0) 出力手段 2 1 1 0 4 は、バッファ中の全データに対応する全評 定値を出力する。ステップ S 2 8 0 9 に行く。

以下、本実施の形態における音声処理装置の具体的な動作について説明する。本実施の 形態において、音素の置換の検知を行う処理が実施の形態 4 等とは異なる。そこで、その 10

20

30

40

異なる処理を中心に説明する。

[0177]

まず、学習者(評価対象者)が、語学学習の開始の指示である動作開始指示を入力する。そして、音声処理装置は、当該動作開始指示を受け付け、次に、例えば、「"う"と発声してください。」と画面出力する。

そして、評価対象者は、"う"と発声し、音声処理装置は、当該発声から、第二ンプリング周波数「32.1 K H z 」を得る。かかる処理は、実施の形態1等において説明した処理と同様である。

[0178]

次に、発声催促部1109は、例えば、「"right"と発声してください。」と画面出力する。そして、学習者は、学習対象の音声「right」を発音する。そして、音声受付部103は、学習者が発音した音声の入力を受け付ける。

次に、サンプリング部106は、受け付けた音声「right」をサンプリング周波数「22.05KHz」でサンプリング処理する。そして、サンプリング部106は、「right」の第一音声データを得る。

[0179]

次に、声道長正規化処理部109は、「right」の第一音声データを第二サンプリング周波数「32.1 KHz」でリサンプリング処理する。そして、声道長正規化処理部109は、第二音声データを得る。次に、音声処理部2610は、第二音声データを、以下のように処理する。

まず、フレーム区分手段1101は、「right」の第二音声データを、短時間フレームに区分する。

そして、フレーム音声データ取得手段 1 1 0 2 は、フレーム区分手段 1 1 0 1 が区分した音声データを、スペクトル分析し、特徴ベクトル系列「 $O = O_1$, O_2 , ・・・, O_T 」を算出する。

次に、評定手段 2 6 1 0 3 の最適状態決定手段 1 1 0 3 1 は、取得した特徴ベクトル系列を構成する各特徴ベクトル o _t に基づいて、所定のフレームの最適状態(特徴ベクトル o _t に対する最適状態)を決定する。

次に、最適状態確率値取得手段11032は、上述した数式1、2により、最適状態における確率値を算出する。

[0180]

次に、評定手段26103は、例えば、最適状態決定手段11031が決定した最適状態を有する音韻全体の状態における1以上の確率値を取得し、当該1以上の確率値の総和をパラメータとして音声の評定値を算出する。つまり、評定手段26103は、例えば、DAPスコアをフレーム毎に算出する。ここで、算出するスコアは、上述したt-p-DAPスコア等でも良い。

[0181]

そして、特殊音声検知手段26101は、算出されたフレームに対応する評定値を用いて、特殊な音声が入力されたか否かを判断する。つまり、評定値(例えば、DAPスコア)が、所定の値より低い区間が存在するか否かを判断する。

[0182]

次に、特殊音声検知手段 2 6 1 0 1 は、図 2 9 に示すように、評定値 (例えば、DAPスコア)が、所定の値より低い区間が、一つの音素内(ここでは音素 2)であるか否かを判断する。そして、一つの音素内で評定値が低ければ、次に、特殊音声検知手段 2 6 1 0 1 は、直前の音素(音素 1)および / または直後の音素(音素 3)に対する評定値(例えば、DAPスコア)を算出し、当該評定値が所定の値より高ければ、音素の置換が発生している可能性があると判断する。次に、特殊音声検知手段 2 6 1 0 1 は、予め格納されている音韻 H M M が一つ存在すれば、音素の置換が発生していると判断する。なお、図 2 9 において、音素 2 において、音素の置換が発生した区間である。なお、図 2 9 において縦軸

10

20

30

40

は評定値であり、当該評定値は、DAP、t-p-DAP等、問わない。

[0183]

次に、評定手段 2 6 1 0 3 は、音素の置換があった旨を示す評定結果(例えば、「音素の置換が発生しました。」)を構成する。そして、出力手段 2 1 1 0 4 は、構成した評定結果を出力する。なお、出力手段 2 1 1 0 4 は、通常の入力音声に対しては、上述したように評定値を出力することが好適である。

[0 1 8 4]

以上、本実施の形態によれば、ユーザが入力した発音を、教師データに対して、如何に似ているかを示す類似度(評定値)を算出し、出力できる。また、かかる場合、本実施の形態によれば、個人差、特に声道長の違いに影響を受けない、精度の高い評定ができる。さらに、本音声処理装置は、特殊音声、特に、音素の置換を検知できるので、極めて精度の高い評定結果が得られる。

[0185]

なお、本実施の形態において、音素の置換を検知できれば良く、評定値の算出アルゴリズムは問わない。評定値の算出アルゴリズムは、上述したアルゴリズム(DAP、t-p-DAP)でも良く、または、本明細書では述べていない他のアルゴリズムでも良い。

[0186]

また、本実施の形態において、音素の置換の検知アルゴリズムは、他のアルゴリズムでも良い。例えば、音素の置換の検知において、所定以上の長さの区間を有することを置換区間の検知で必須としても良い。その他、置換の検知アルゴリズムの詳細は種々考えられる。

[0187]

さらに、本実施の形態における音声処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータに、第一サンプリング周波数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリングステップと、第二サンプリング周波数「第一サンプリング周波数/(教師データフォルマント周波数/評価対象者フォルマント周波数)」で、前記音声受付ステップで受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理ステップと、前記第二音声データを処理する音声処理ステップを実行させるためのプログラム、である

[0188]

また、上記プログラムにおいて、音声処理ステップは、前記第二音声データを、フレームに区分するフレーム区分ステップと、前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム音声データ取得ステップと、前記フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する特殊音声検知ステップと、教師データと前記入力音声データと前記特殊音声検知ステップにおける検知結果に基づいて、前記受け付けた音声の評定を行う評定ステップと、前記評定ステップにおける評定結果を出力する出力ステップを具備する、ことは好適である。

[0189]

また、上記プログラムにおいて、特殊音声検知ステップは、一の音素の評定値が所定の条件を満たすことを検知し、特殊音声検知ステップで前記所定の条件を満たすことを検知した場合に、少なくとも音素の置換があった旨を示す評定結果を構成する、ことは好適である。

(実施の形態6)

[0190]

本実施の形態において、入力音声において、特殊音声を検知し、比較対象の音声と入力音声の類似度を精度高く評定できる音声処理装置について説明する。特に、本音声処理装置は、音韻の欠落を検知できる音声処理装置である。

[0191]

図30は、本実施の形態における音声処理装置のブロック図である。本音声処理装置は

10

20

30

50

10

20

30

40

50

、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部104、第一サンプリング周波数格納部105、サンプリング部106、評価対象者フォルマント周波数取得部107、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部3010、発声催促部1109を具備する

[0192]

音声処理部3010は、フレーム区分手段1101、フレーム音声データ取得手段1102、特殊音声検知手段30101、評定手段30103、出力手段21104を具備する。なお、評定手段30103は、最適状態決定手段11031、最適状態確率値取得手段11032を具備する。

[0193]

特殊音声検知手段30101は、一の音素の評定値が所定の値より低く、かつ当該音素の直前の音素または当該音素の直後の音素の評定値が所定の値より高いことを検知する。また、特殊音声検知手段30101は、一の音素の評定値が所定の値より低く、かつ当該音素の区間長が所定の長さよりも短いことを検知しても良い。また、特殊音声検知手段30101は、直前の音素に対応する確率値、または直後の音素に対応する確率値が、当該000音素の確率値より高いことを検知しても良い。かかる場合に、特殊音声検知手段30101は、音韻の欠落を検知することは好適である。さらに、音素の区間長が所定の長さりも短いことを欠落の条件に含めることにより、音韻の欠落の検知の精度は向上する。特殊音声検知手段30101の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0194]

評定手段30103は、特殊音声検知手段30101が所定の条件を満たすことを検知した場合に、少なくとも音素の欠落があった旨を示す評定結果を構成する。評定手段30103は、通常、MPUやメモリ等から実現され得る。評定手段30103の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0195]

次に、音声処理装置の動作について、図31、図32のフローチャートを用いて説明する。なお、図31のフローチャートは、図12のフローチャートと比較して、ステップS3101の評定処理のみが異なるので、図31のフローチャートの説明は省略する。ステップS3101の評定処理の詳細について、図32のフローチャートを用いて説明する。図32のフローチャートにおいて、図2、図19、図23、図28のフローチャートの処理と同様の処理については、その説明を省略する。

[0196]

(ステップS3201)特殊音声検知手段30101は、バッファに蓄積されているデータに対して、直前の音素に対応する教師データの確率値または、直後の音素に対応する教師データの確率値が、予定されている音素に対応する教師データの確率値より高いか否かを判断する。高ければステップS3202に行き、高くなければステップS2810に行く。なお、ステップS3202に行くための条件として、バッファに蓄積されているデータに対応するフレーム音声データ群の区間長が所定の長さ以下であることを付加しても良い。

(ステップS3202) 評定手段30103は、音素の欠落があった旨を示す評定結果 を構成する。ステップS2808に行く。

なお、図32のフローチャートにおいて、評定対象の音素(欠落したであろう音素)の 区間長が、所定の長さ(例えば、3フレーム)よりも短いことを条件としても良いし、か かる条件は無くても良い。

以下、本実施の形態における音声処理装置の具体的な動作について説明する。本実施の

形態において、音素の欠落の検知を行う処理が実施の形態 5 等とは異なる。そこで、その 異なる処理を中心に説明する。

[0197]

まず、学習者(評価対象者)が、語学学習の開始の指示である動作開始指示を入力する。そして、音声処理装置は、当該動作開始指示を受け付け、次に、例えば、「"う"と発声してください。」と画面出力する。

そして、評価対象者は、"う"と発声し、音声処理装置は、当該発声から、第二ンプリング周波数「32.1 K H z 」を得る。かかる処理は、実施の形態1等において説明した処理と同様である。

[0198]

次に、発声催促部 1 1 0 9 は、例えば、「"right"と発声してください。」と画面出力する。そして、学習者は、学習対象の音声「right」を発音する。そして、音声受付部 1 0 3 は、学習者が発音した音声の入力を受け付ける。

次に、サンプリング部106は、受け付けた音声「right」をサンプリング周波数「22.05KHz」でサンプリング処理する。そして、サンプリング部106は、「right」の第一音声データを得る。

[0199]

次に、声道長正規化処理部109は、「right」の第一音声データを第二サンプリング周波数「32.1 KHz」でリサンプリング処理する。そして、声道長正規化処理部109は、第二音声データを得る。次に、音声処理部3010は、第二音声データを、以下のように処理する。

まず、フレーム区分手段1101は、「right」の第二音声データを、短時間フレームに区分する。

そして、フレーム音声データ取得手段 1 1 0 2 は、フレーム区分手段 1 1 0 1 が区分した音声データを、スペクトル分析し、特徴ベクトル系列「 $O = O_1$, O_2 , ・・・, O_T 」を算出する。

次に、最適状態決定手段 1 1 0 3 1 は、取得した特徴ベクトル系列を構成する各特徴ベクトル o_t に基づいて、所定のフレームの最適状態(特徴ベクトル o_t に対する最適状態)を決定する。

次に、最適状態確率値取得手段11032は、上述した数式1、2により、最適状態における確率値を算出する。

[0200]

次に、評定手段30103は、例えば、最適状態決定手段11031が決定した最適状態を有する音韻全体の状態における1以上の確率値を取得し、当該1以上の確率値の総和をパラメータとして音声の評定値を算出する。つまり、評定手段30103は、例えば、DAPスコアをフレーム毎に算出する。ここで、算出するスコアは、上述したt-p-DAPスコア等でも良い。

[0201]

そして、特殊音声検知手段30101は、算出されたフレームに対応する評定値を用いて、特殊な音声が入力されたか否かを判断する。つまり、評定値(例えば、DAPスコア)が、所定の値より低い区間が存在するか否かを判断する。

[0202]

次に、特殊音声検知手段30101は、図33に示すように、評定値(例えば、DAP スコア)が、所定の値より低い区間が、一つの音素内(ここでは音素2)であるか否かを 判断する。そして、一つの音素内で評定値が低ければ、特殊音声検知手段30101は、 直前の音素(音素1)または直後の音素(音素3)に対する評定値(例えば、DAPスコ ア)を算出し、当該評定値が所定の値より高ければ、音素の欠落が発生している可能性が あると判断する。そして、当該区間長が、例えば、3フレーム以下の長さであれば、かか る音素は欠落したと判断する。なお、図33において、音素2の欠落が発生したことを示 10

20

30

40

す。なお、図33において縦軸は評定値であり、当該評定値は、DAP、t-p-DAP 等、問わない。また、上記区間長の所定値は、「3フレーム以下」ではなく、「5フレー ム以下」でも、「6フレーム以下」でも良い。

[0203]

次に、評定手段30103は、音素の欠落があった旨を示す評定結果(例えば、「音素の欠落が発生しました。」)を構成する。そして、出力手段21104は、構成した評定結果を出力する。なお、出力手段21104は、通常の入力音声に対しては、上述したように評定値を出力することが好適である。

[0204]

以上、本実施の形態によれば、ユーザが入力した発音を、教師データに対して、如何に似ているかを示す類似度(評定値)を算出し、出力できる。また、かかる場合、本実施の形態によれば、個人差、特に声道長の違いに影響を受けない、精度の高い評定ができる。さらに、本音声処理装置は、特殊音声、特に、音素の欠落を検知できるので、極めて精度の高い評定結果が得られる。

[0205]

なお、本実施の形態において、音素の欠落を検知できれば良く、評定値の算出アルゴリズムは問わない。評定値の算出アルゴリズムは、上述したアルゴリズム(DAP、t-p-DAP)でも良く、または、本明細書では述べていない他のアルゴリズムでも良い。

[0206]

また、本実施の形態において、音素の欠落の検知アルゴリズムは、他のアルゴリズムでも良い。例えば、音素の欠落の検知において、所定長さ未満の区間であることを欠落区間の検知で必須としても良いし、区間長を考慮しなくても良い。

[0207]

さらに、本実施の形態における音声処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータに、第一サンプリング周波数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリングステップと、第二サンプリング周波数「第一サンプリング周波数 / (教師データフォルマント周波数 / 評価対象者フォルマント周波数)」で、前記音声受付ステップで受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理ステップと、前記第二音声データを処理する音声処理ステップを実行させるためのプログラム、である

[0208]

また、上記プログラムにおいて、音声処理ステップは、前記第二音声データを、フレームに区分するフレーム区分ステップと、前記区分されたフレーム毎の音声データであるフレーム音声データを1以上得るフレーム音声データ取得ステップと、前記フレーム毎の入力音声データに基づいて、特殊な音声が入力されたことを検知する特殊音声検知ステップと、教師データと前記入力音声データと前記特殊音声検知ステップにおける検知結果に基づいて、前記受け付けた音声の評定を行う評定ステップと、前記評定ステップにおける評定結果を出力する出力ステップを具備する、ことは好適である。

[0209]

また、上記プログラムにおいて、特殊音声検知ステップは、一の音素の評定値が所定の 条件を満たすことを検知し、特殊音声検知ステップで前記所定の条件を満たすことを検知 した場合に、少なくとも音素の欠落があった旨を示す評定結果を構成する、ことは好適で ある。

(実施の形態7)

[0210]

本実施の形態における音声処理装置は、サンプリング周波数を変更し、リサンプリングを行わずに評定した場合の評定値と、リサンプリングを行って評定した場合の評定値とを取得し、2つの評定値に基づいて、最終的な評定値を算出する音声処理装置である。例えば、本音声処理装置は、2つの評定値の平均値を最終的な評定値としても良いし、2つの

10

20

30

40

評定値の最大値を最終的な評定値としても良い。また、本音声処理装置は、例えば、カラオケ評定装置である。

[0211]

図34は、本実施の形態における音声処理装置のブロック図である。本音声処理装置は、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部105、サンプリング部106、評価対象者フォルマント周波数格納部107、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部3410、発声催促部1109を具備する

音声処理部3410は、フレーム区分手段34101、フレーム音声データ取得手段34102、評定手段34103、出力手段1104を具備する。

評定手段34103は、第一評定手段341031、第二評定手段341032、評定結果取得手段341033を具備する。

フレーム区分手段 3 4 1 0 1 は、音声をフレームに区分し、かつ、前記第二音声データをフレームに区分する。

[0212]

フレーム音声データ取得手段 3 4 1 0 2 は、音声が区分されたフレーム毎の音声データである第一フレーム音声データを1以上得て、かつ前記第二音声データが区分されたフレーム毎の音声データである第二フレーム音声データを1以上得る。

[0213]

評定手段34103は、教師データと1以上のフレーム音声データに基づいて、音声受付部103が受け付けた音声の評定を行う。評定手段34103は、以下の第一評定手段341031の評定結果と、第二評定手段341032の評定結果に基づいて、最終的な評定結果を得る。

第一評定手段341031は、教師データと1以上の第一フレーム音声データに基づいて、音声受付部が受け付けた音声の評定を行う。

第二評定手段341032は、教師データと1以上の第二フレーム音声データに基づいて、音声受付部が受け付けた音声の評定を行う。

[0214]

評定結果取得手段 3 4 1 0 3 3 は、第一評定手段 3 4 1 0 3 1 における評定結果(以下、適宜「第一評定結果」という。)と第二評定手段 3 4 1 0 3 2 における評定結果(以下、適宜「第二評定結果」という。)に基づいて、最終的な評定結果を得る。評定結果取得手段 3 4 1 0 3 3 は、例えば、第一評定結果と第二評定結果の平均値を、最終的な評定結果としても良いし、第一評定結果と第二評定結果の大きい方の値を最終的な評定結果としても良いし、第一評定結果と第二評定結果の小さい方の値を最終的な評定結果としても良い。

[0215]

フレーム区分手段34101、フレーム音声データ取得手段34102、第一評定手段341031、第二評定手段341032、評定結果取得手段341033は、通常、MPUやメモリ等から実現され得る。フレーム区分手段34101等の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

次に、音声処理装置の動作について図35のフローチャートを用いて説明する。図35のフローチャートにおいて、図2、図12のフローチャートと異なるステップについてのみ説明する。

[0216]

(ステップS 3 5 0 1) 第一評定手段 3 4 1 0 3 1 は、第一評定処理を行う。第一評定処理とは、教師データと1以上の第一フレーム音声データに基づいて、音声受付部 1 0 3 が受け付けた音声の評定を行う処理である。第一評定処理は、リサンプリングしない第一音声データを評定する処理である。第一評定処理における評定のアルゴリズムは、上記の

10

20

30

40

実施の形態 1 から実施の形態 6 で述べたいずれのアルゴリズム(DAP、t-p-DAP、無音区間考慮、挿入を考慮、置換を考慮、欠落を考慮など)または、それらを組み合わせたアルゴリズムでも良い。

[0217]

(ステップS3502)第二評定手段341032は、第二評定処理を行う。第二評定処理とは、教師データと1以上の第二フレーム音声データに基づいて、音声受付部103が受け付けた音声の評定を行う処理である。第二評定処理は、リサンプリングした第二音声データを評定する処理である。第二評定処理における評定のアルゴリズムは、上記の実施の形態1から実施の形態6で述べたいずれのアルゴリズム(DAP、t-p-DAP、無音区間考慮、挿入を考慮、置換を考慮、欠落を考慮など)または、それらを組み合わせたアルゴリズムでも良い。なお、第一評定処理と第二評定処理のアルゴリズムは、同一であることが好適である。

[0218]

(ステップS3503)評定結果取得手段341033は、第一評定手段341031における評定結果(第一評定結果)と第二評定手段341032における評定結果(第二評定結果)に基づいて、最終的な評定結果を得る。評定結果取得手段341033は、例えば、第一評定結果と第二評定結果の評定値のうち、高得点の方の評定値を最終的な評定結果とする。

[0219]

以上、本実施の形態によれば、ユーザが入力した発音を、教師データに対して、如何に似ているかを示す類似度(評定値)を算出し、出力できる。また、かかる場合、本実施の形態によれば、個人差を考慮した精度の高い評定ができる。さらに、本音声処理装置は、個人差を考慮した評定と、個人差を考慮しない評定の両方を利用した評定が行える。つまり、本実施の形態によれば、例えば、第一評定結果と第二評定結果の評定値のうち、高得点の方の評定値を最終的な評定結果とすることができ、カラオケ評定装置等として有効である。

(実施の形態8)

本実施の形態における音声処理装置の音声処理は、音声認識である。

図36は、本実施の形態における音声処理装置のブロック図である。

[0220]

本音声処理装置は、入力受付部101、教師データ格納部102、音声受付部103、教師データフォルマント周波数格納部104、第一サンプリング周波数格納部105、サンプリング部106、評価対象者フォルマント周波数取得部107、評価対象者フォルマント周波数格納部108、声道長正規化処理部109、音声処理部3610、発声催促部1109を具備する。

音声処理部3610は、音声認識手段36101、出力手段36102を具備する。

[0221]

音声処理部3610の音声認識手段36101は、第二音声データに基づいて、音声認識処理を行う。音声認識のアルゴリズムは、問わない。音声認識処理は、公知のアルゴリズムで良い。本実施の形態において、リサンプリングした第二音声データに基づいて音声認識することにより、精度の高い音声認識が可能である。音声処理部3610は、通常、MPUやメモリ等から実現され得る。音声処理部3610の処理手順は、通常、ソフトウェアで実現され、当該ソフトウェアはROM等の記録媒体に記録されている。但し、ハードウェア(専用回路)で実現しても良い。

[0222]

出力手段36102は、音声認識結果を出力する。ここで、出力とは、ディスプレイへの表示、プリンタへの印字、音出力、外部の装置への送信、記録媒体への蓄積等を含む概念である。出力手段36102は、ディスプレイやスピーカー等の出力デバイスを含むと考えても含まないと考えても良い。出力手段36102は、出力デバイスのドライバーソフトまたは、出力デバイスのドライバーソフトと出力デバイス等で実現され得る。

10

20

30

40

次に、音声処理装置の動作について図37のフローチャートを用いて説明する。なお、図37のフローチャートにおいて、図2、図12のフローチャートの処理と同様の処理については、その説明を省略する。

[0223]

(ステップS3701)音声認識手段36101は、ステップS1208でリサンプリング処理され、得られた第二音声データに基づいて、音声認識処理を行う。なお、音声認識手段36101は、教師データとのマッチングを取り、教師データに近い音であると認識することにより、認識結果を得る。

(ステップ S 3 7 0 2) 出力手段 3 6 1 0 2 は、ステップ S 3 7 0 1 における音声認識 結果を出力する。ステップ S 1 2 0 6 に戻る。

以上、本実施の形態によれば、精度高く音声認識できる。

[0224]

なお、本実施の形態における音声処理装置を実現するソフトウェアは、以下のようなプログラムである。つまり、このプログラムは、コンピュータに、第一サンプリング周波数で、受け付けた音声をサンプリングし、第一音声データを取得するサンプリングステップと、第二サンプリング周波数「第一サンプリング周波数/(教師データフォルマント周波数/評価対象者フォルマント周波数)」で、前記音声受付ステップで受け付けた音声に対して、サンプリング処理を行い、第二音声データを得る声道長正規化処理ステップと、前記第二音声データに基づいて、音声認識処理を行う音声処理ステップを実行させるためのプログラム、である。

[0225]

また、上記の実施の形態において検出した特殊音声は、無音、挿入、置換、欠落であった。音声処理装置は、かかるすべての特殊音声について検知しても良いことはいうまでもない。また、音声処理装置は、主として、実施の形態 1、実施の形態 2 において述べた評定値の算出アルゴリズムを利用して、特殊音声の検出を行ったが、他の評定値の算出アルゴリズムを利用しても良い。

[0226]

また、特殊音声は、無音、挿入、置換、欠落に限られない。例えば、特殊音声は、garbage(雑音などの雑多な音素等)であっても良い。受け付けた音声にgarbageが混入している場合、その区間は類似度の計算対象から除外するのがしばしば望ましい。例えば、発音評定においては、学習者の発声には通常、息継ぎや無声区間などが数多く表れ、それらに対応する発声区間を評定対象から取り除くことが好適である。なお、無音は、一般に、garbageの一種である、と考える。

[0227]

そこで、どの音素にも属さない雑多な音素(garbage音素)を設定し、garbageのHMMをあらかじめ格納しておく。スコア低下区間において、garbageのHMMに対する評定値(t(j))が所定の値より大きい場合、その区間はgarbage区間と判定することは好適である。特に、発音評定において、garbage区間が2つの単語にまたがっている場合、息継ぎなどが起こったものとして、評定値の計算対象から除外することは極めて好適である。

[0228]

また、図38は、本明細書で述べたプログラムを実行して、上述した種々の実施の形態の音声処理装置を実現するコンピュータの外観を示す。上述の実施の形態は、コンピュータハードウェア及びその上で実行されるコンピュータプログラムで実現され得る。図38は、このコンピュータシステム340の概観図であり、図39は、コンピュータシステム340のブロック図である。

[0229]

図 38 において、コンピュータシステム 340 は、 FD (Flexible Disk) ドライブ、 CD - ROM (Compact Disk Read Only Memory) ドライブを含むコンピュータ 341 と、キーボード 342 と、マウス 343 と、モ

10

20

30

40

ニタ344と、マイク345とを含む。

[0230]

図39において、コンピュータ341は、FDドライブ3411、CD-ROMドライブ3412に加えて、CPU(Central Processing Unit)3413と、CPU3413、CD-ROMドライブ3412及びFDドライブ3411に接続されたバス3414と、ブートアッププログラム等のプログラムを記憶するためのROM(Read-Only Memory)3415と、CPU3413に接続され、アプリケーションプログラムの命令を一時的に記憶するとともに一時記憶空間を提供するためのRAM(Random Access Memory)3416と、アプリケーションプログラム、システムプログラム、及びデータを記憶するためのハードディスク3417とを含む。ここでは、図示しないが、コンピュータ341は、さらに、LANへの接続を提供するネットワークカードを含んでも良い。

[0231]

コンピュータシステム 3 4 0 に、上述した実施の形態の音声処理装置の機能を実行させるプログラムは、CD-ROM 3 5 0 1、またはFD 3 5 0 2 に記憶されて、CD-ROMドライブ 3 4 1 2 またはFDドライブ 3 4 1 1 に挿入され、さらにハードディスク 3 4 1 7 に転送されても良い。これに代えて、プログラムは、図示しないネットワークを介してコンピュータ 3 4 1 に送信され、ハードディスク 3 4 1 7 に記憶されても良い。プログラムは実行の際にRAM 3 4 1 6 にロードされる。プログラムは、CD-ROM 3 5 0 1、FD 3 5 0 2 またはネットワークから直接、ロードされても良い。

[0232]

プログラムは、コンピュータ341に、上述した実施の形態の音声処理装置の機能を実行させるオペレーティングシステム(OS)、またはサードパーティープログラム等は、必ずしも含まなくても良い。プログラムは、制御された態様で適切な機能(モジュール)を呼び出し、所望の結果が得られるようにする命令の部分のみを含んでいれば良い。コンピュータシステム340がどのように動作するかは周知であり、詳細な説明は省略する。

[0233]

また、上記各実施の形態において、各処理(各機能)は、単一の装置(システム)によって集中処理されることによって実現されてもよく、あるいは、複数の装置によって分散処理されることによって実現されてもよい。

また、上記のプログラムを実行するコンピュータは、単数であってもよく、複数であってもよい。すなわち、集中処理を行ってもよく、あるいは分散処理を行ってもよい。

本発明は、以上の実施の形態に限定されることなく、種々の変更が可能であり、それらも本発明の範囲内に包含されるものであることは言うまでもない。

【産業上の利用可能性】

[0234]

以上のように、本発明にかかる音声処理装置は、評価対象者の話者特性に応じた精度の高い音声処理ができるという効果を有し、発音評定装置やカラオケ評定装置や音声認識装置等として有用である。

【図面の簡単な説明】

[0235]

- 【図1】実施の形態1における音声処理装置のブロック図
- 【図2】同音声処理装置の動作について説明するフローチャート
- 【図3】同声道長正規化処理について説明するフローチャート
- 【図4】同HMMの仕様の例を示す図
- 【図5】同F1、F2の計測結果を示す図
- 【図6】同音声分析条件を示す図
- 【図7】同算出した評定値をグラフで表した例を示す図
- 【図8】同算出した評定値をグラフで表した例を示す図
- 【図9】同出力例を示す図

20

10

30

40

50

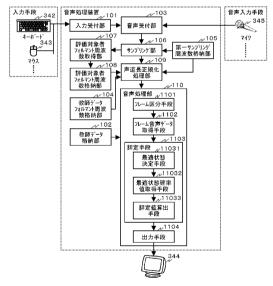
```
【図10】同出力例を示す図
【図11】実施の形態2における音声処理装置のブロック図
【図12】同音声処理装置の動作について説明するフローチャート
【図13】同第二サンプリング周波数算出処理について説明するフローチャート
【図14】同評定処理について説明するフローチャート
【図15】同評定結果(t-p-DAPスコア)を示す図
【図16】同出力例を示す図
【図17】実施の形態3における音声処理装置のブロック図
【図18】同音声処理装置の動作について説明するフローチャート
                                               10
【図19】同評定処理について説明するフローチャート
【図20】同無音データの検知について説明する図
【図21】実施の形態4における音声処理装置のブロック図
【図22】同音声処理装置の動作について説明するフローチャート
【図23】同評定処理について説明するフローチャート
【図24】同音素の挿入の検知について説明する図
【図25】同出力例を示す図
【図26】実施の形態5における音声処理装置のブロック図
【図27】同音声処理装置の動作について説明するフローチャート
【図28】同評定処理について説明するフローチャート
                                               20
【図29】同音素の置換の検知について説明する図
【図30】実施の形態6における音声処理装置のブロック図
【図31】同音声処理装置の動作について説明するフローチャート
【図32】同評定処理について説明するフローチャート
【図33】同音素の欠落の検知について説明する図
【図34】実施の形態7における音声処理装置のブロック図
【図35】同音声処理装置の動作について説明するフローチャート
【図36】実施の形態8における音声処理装置のブロック図
【図37】同音声処理装置の動作について説明するフローチャート
【図38】同音声処理装置を構成するコンピュータシステムの概観図
                                               30
【図39】同音声処理装置を構成するコンピュータのブロック図
【符号の説明】
[0236]
 1 0 1
     入力受付部
 102 教師データ格納部
 103 音声受付部
     教師データフォルマント周波数格納部
 1 0 4
 1 0 5 第一サンプリング周波数格納部
 1 0 6
     サンプリング部
     評価対象者フォルマント周波数取得部
 1 0 7
                                               40
 108 評価対象者フォルマント周波数格納部
 109 声道長正規化処理部
 110, 1110, 1710, 2110, 2610, 3010, 3410, 3610
音声処理部
 1 1 0 1 、 3 4 1 0 1 フレーム区分手段
 1 1 0 2 、 3 4 1 0 2 フレーム音声データ取得手段
 1 1 0 3 、 1 1 1 0 3 、 1 7 1 0 3 、 2 1 1 0 3 、 2 6 1 0 3 、 3 0 1 0 3 、 3 4 1 0
3 評定手段
 1 1 0 4 、 2 1 1 0 4 、 3 6 1 0 2
                   出力手段
 1 1 0 9 発声催促部
```

11031 最適状態決定手段

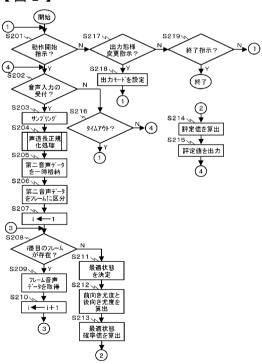
- 11032 最適状態確率値取得手段
- 17101、21101、26101、30101 特殊音声検知手段
- 36101 音声認識手段
- 111032 発音区間フレーム音韻確率値取得手段
- 171011 無音データ格納手段
- 171012 無音区間検出手段
- 3 4 1 0 3 1 第一評定手段
- 3 4 1 0 3 2 第二評定手段
- 3 4 1 0 3 3 評定結果取得手段

10

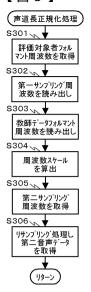
【図1】



【図2】



【図3】



【図4】

HMM\$17°	混合がウス型, mono phone (対角共分散行列)		
状態数(モテ゚ル毎)	3		
混合数(状態毎)	8		
HMM総数	44		
総状態数	132		

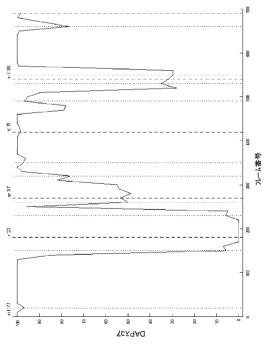
【図5】

	F1(Hz)	F2(Hz)
成人男性(139名の平均)	340	1184
成人女性(92名の平均)	370	1310
少年12-18歳(69名の平均)	402	1534
子供5-11歳(56名の平均)	510	1725

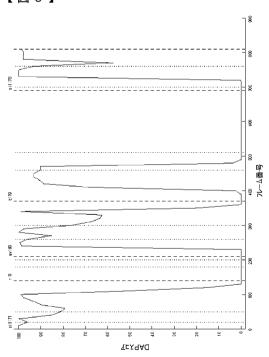
【図6】

プリエンファシス	1-0.97 _z -1		
窓関数	Hamming窓		
分析フレーム長	25msec		
フレーム周期	10msec		
特徴パラメータ	MFCC(39次元)		
	L		

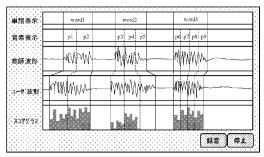
【図7】



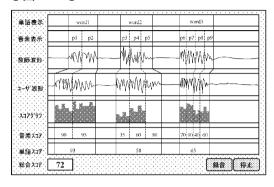
【図8】



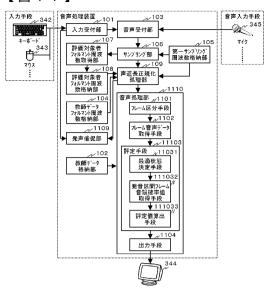
【図9】



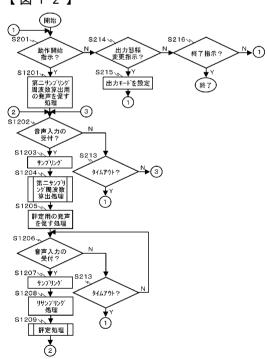
【図10】



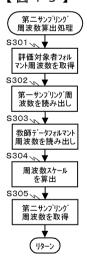
【図11】

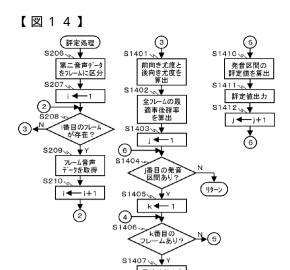


【図12】



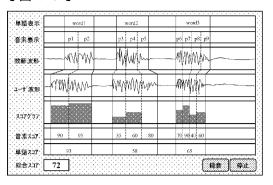
【図13】





最適状態を含む音韻の全状態の確率値を 取得

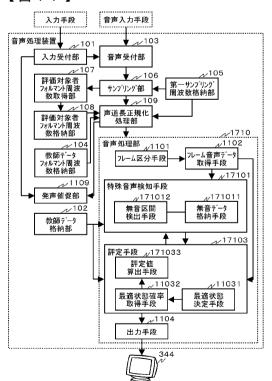
【図16】



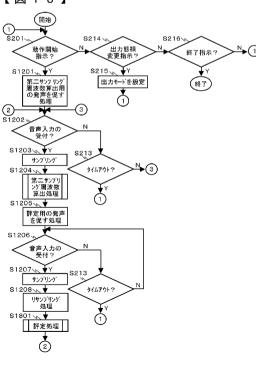
【図15】

	Phoneme /r/ /ay/ /t/			Word /right/
アメリカ人男性	39	99	100	84
日本人男性	0	70	22	33

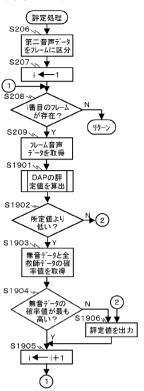


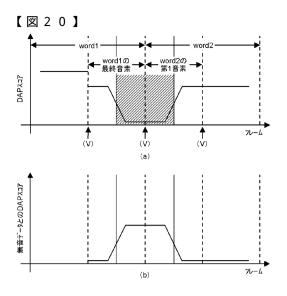


【図18】

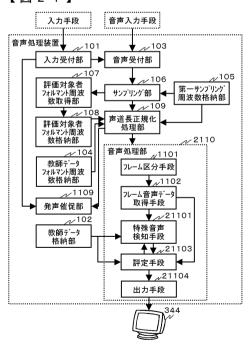


【図19】

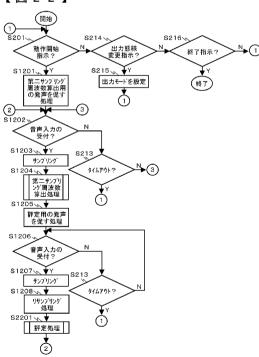


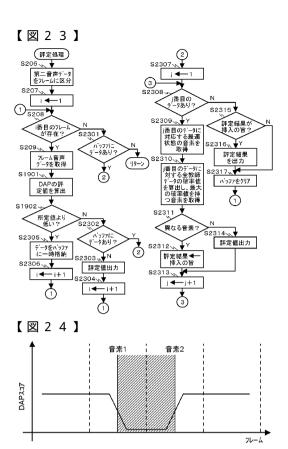


【図21】

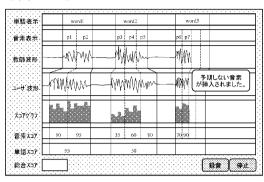


【図22】

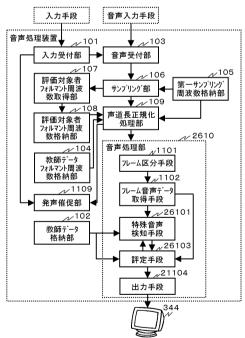




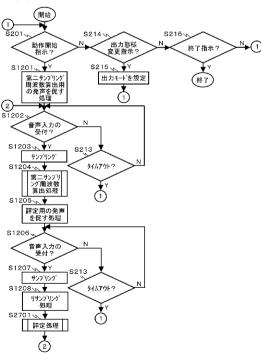
【図25】



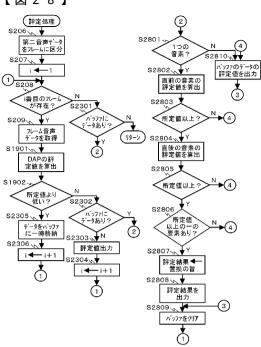




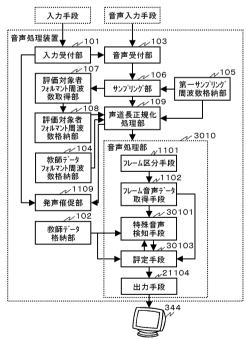
【図27】 開始

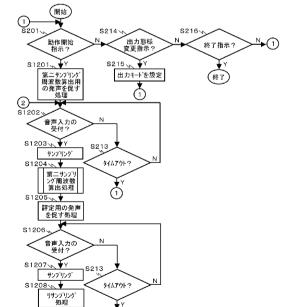


【図28】



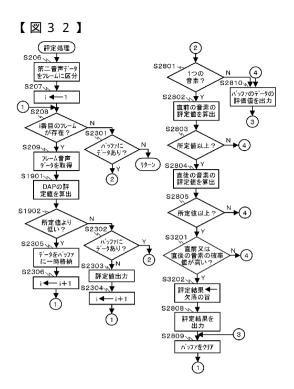
【図30】

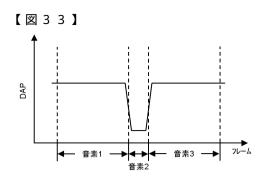


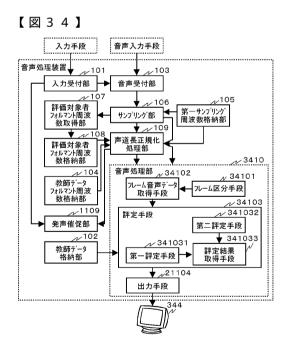


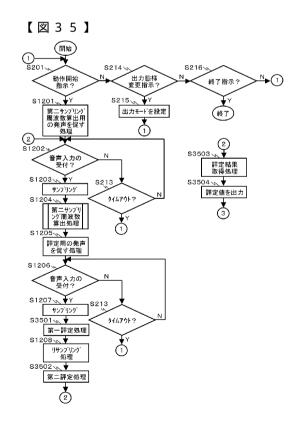
Ť

【図31】

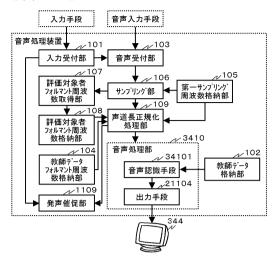


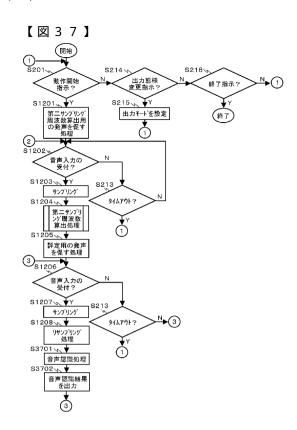




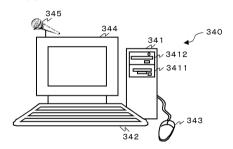


【図36】

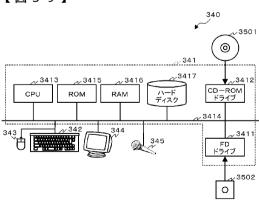




【図38】



【図39】



フロントページの続き

(72)発明者 加藤 宏明

京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内

審査官 毛利 太郎

(56)参考文献 特開2001-042889(JP,A)

特開昭62-174798(JP,A)

特表2002-515136(JP,A)

特開2001-265211(JP,A)

特開平10-222190(JP,A)

特開平06-110494(JP,A)

特開2006-227030(JP,A)

特開2001-117598(JP,A)

特開平11-259081(JP,A)

(58)調査した分野(Int.CI., DB名)

G10L 11/00-21/06