

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4775788号  
(P4775788)

(45) 発行日 平成23年9月21日(2011.9.21)

(24) 登録日 平成23年7月8日(2011.7.8)

(51) Int. Cl.		F I	
G 0 9 B	19/00 (2006.01)	G 0 9 B	19/00 G
G 0 9 B	19/04 (2006.01)	G 0 9 B	19/04
G 1 0 L	15/00 (2006.01)	G 1 0 L	15/00 2 0 0 E
G 1 0 L	15/04 (2006.01)	G 1 0 L	15/04 3 0 0 C
G 1 0 L	15/14 (2006.01)	G 1 0 L	15/14 2 0 0 Z

請求項の数 4 (全 41 頁)

(21) 出願番号 特願2005-16424 (P2005-16424)  
 (22) 出願日 平成17年1月25日(2005.1.25)  
 (65) 公開番号 特開2006-227030 (P2006-227030A)  
 (43) 公開日 平成18年8月31日(2006.8.31)  
 審査請求日 平成19年8月3日(2007.8.3)  
 (31) 優先権主張番号 特願2005-13157 (P2005-13157)  
 (32) 優先日 平成17年1月20日(2005.1.20)  
 (33) 優先権主張国 日本国(JP)

(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 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  
 (72) 発明者 渡辺 秀行  
 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 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

最終頁に続く

(54) 【発明の名称】 発音評定装置、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、フレーム毎に状態識別子と状態間を遷移することで得られる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手段と、

音声の入力を受け付ける音声受付部と、

前記音声受付部が受け付けた音声を、フレームに区分するフレーム区分部と、

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得るフレーム音声データ取得部と、

前記フレーム毎のフレーム音声データに基づいて、音素の挿入を検知する特殊音声検知部と、

前記教師データと前記フレーム毎のフレーム音声データと前記特殊音声検知部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評定部と、

前記評定部の評定結果を出力する出力部を具備し、

前記特殊音声検知部は、

前記フレーム音声データ取得部が得た複数のフレーム音声データから、1以上の音素を取得し、当該1以上のいずれかの音素について、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値より低い場合、または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、音素の挿入を検知し、

10

20

前記評定部は、

1以上のフレーム音声データに対する、前記教師データの状態間を遷移することで得られる遷移確率の累積が最も高い経路の各状態である1以上の最適状態を決定する最適状態決定部と、

前記最適状態決定部が決定した1以上の最適状態の事後確率を示す確率値を取得する最適状態確率値取得部と、

前記最適状態確率値取得部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部とを具備し、

前記特殊音声検知部が音素の挿入を検知した場合に、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発音評定装置。

10

【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、フレーム毎に状態識別子と状態間を遷移することで得られる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手段と、

音声の入力を受け付ける音声受付部と、

前記音声受付部が受け付けた音声を、フレームに区分するフレーム区分部と、

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得るフレーム音声データ取得部と、

前記フレーム毎のフレーム音声データに基づいて、音素の置換を検知する特殊音声検知部と、

20

前記教師データと前記フレーム毎のフレーム音声データと前記特殊音声検知部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評定部と、

前記評定部の評定結果を出力する出力部を具備し、

前記特殊音声検知部は、

前記フレーム音声データ取得部が得た複数のフレーム音声データから、1以上の音素を取得し、当該1以上のいずれかの音素について、一の音素の評定値が所定値より低く、当該一の音素の直前の音素の評定値または直後の音素の評定値が所定の値より高い場合に、音素の置換を検知し、

前記評定部は、

1以上のフレーム音声データに対する、前記教師データの状態間を遷移することで得られる遷移確率の累積が最も高い経路の各状態である1以上の最適状態を決定する最適状態決定部と、

30

前記最適状態決定部が決定した1以上の最適状態の事後確率を示す確率値を取得する最適状態確率値取得部と、

前記最適状態確率値取得部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部とを具備し、

前記特殊音声検知部が音素の置換を検知した場合に、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発音評定装置。

【請求項3】

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、フレーム毎に状態識別子と状態間を遷移することで得られる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手段と、

40

音声の入力を受け付ける音声受付部と、

前記音声受付部が受け付けた音声を、フレームに区分するフレーム区分部と、

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得るフレーム音声データ取得部と、

前記フレーム毎のフレーム音声データに基づいて、音素の欠落を検知する特殊音声検知部と、

前記教師データと前記フレーム毎のフレーム音声データと前記特殊音声検知部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評定部と、

50

前記評定部の評定結果を出力する出力部を具備し、

前記特殊音声検知部は、

前記フレーム音声データ取得部が得た複数のフレーム音声データから、1以上の音素を取得し、当該1以上のいずれかの音素について、一の音素の評定値が所定値より低く、当該一の音素の直前の音素の評定値または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に、音素の欠落を検知し、

前記評定部は、

1以上のフレーム音声データに対する、前記教師データの状態間を遷移することで得られる遷移確率の累積が最も高い経路の各状態である1以上の最適状態を決定する最適状態決定部と、

10

前記最適状態決定部が決定した1以上の最適状態の事後確率を示す確率値を取得する最適状態確率値取得部と、

前記最適状態確率値取得部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部とを具備し、

前記特殊音声検知部が音素の欠落を検知した場合に、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発音評定装置。

【請求項4】

前記評定値算出部は、

前記最適状態確率値取得部が取得した最適状態の確率値と、当該最適状態の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をパラメータとして音声の評定値を算出する請求項1から請求項3いずれか記載の発音評定装置。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、入力された音声进行评估する装置等に関し、特に、語学学習等に利用できる発音評定装置等に関するものである。

【背景技術】

【0002】

従来技術として、以下の語学学習装置がある（特許文献1参照）。本語学学習装置は、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。本従来技術の語学学習装置は、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。

30

【特許文献1】特開2003-228279（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】

【0003】

しかし、従来技術においては、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。

40

【0004】

一方、一般に、入力される音声には、想定されている音韻列の順序に従わない音韻の挿入や置換、および本来あるはずの音韻の欠落が起こり得る。例えば、習得段階のネイティブではない学習者の発音には、正解の音声のデータにはない虚偽の音韻の挿入、他の音韻への置換、および本来なければならぬ音韻の欠落が起こる。また、学習者の発音には息継ぎなどの無声区間が多数存在する。従来技術においては、スコアが低い場合、受け付けた音声の音韻列は正しいが、単に類似度が低下したのか、音韻の挿入や置換あるいは「garbage」（どのカテゴリにも含まれない雑多な音韻）により低下したのかが判断

50

できない。さらに、従来の技術においては、ごく短区間のみスコアが低下している場合、類似度が低いのか、正しい音韻が欠落しているのか、あるいは *g a r b a g e* が混入しているのかが判断できない。

【課題を解決するための手段】

【0005】

本第一の発明の発音評定装置は、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、音声の入力を受け付ける音声受付部と、前記音声受付部で受け付けた音声を、フレームに区分するフレーム区分部と、前記区分されたフレーム毎の入力音声データを得る入力音声データ取得部と、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部と、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部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評定部と、前記評定部の評定結果を出力する出力部を具備する発音評定装置である。

10

かかる発音評定装置は、特殊な音声が入力されたことを検知でき、当該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える。

【0006】

また、本第二の発明の発音評定装置は、第一の発明の発音評定装置において、特殊音声検知部は、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発音評定装置である。

20

かかる構成により、無音の区間を検出でき、精度が高い音声の評定が可能となる。

【0007】

また、本第三の発明の発音評定装置は、第一の発明の発音評定装置において、特殊音声検知部は、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、前記評定部は、前記特殊音声検知部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発音評定装置である。

かかる構成により、音素の挿入を検出でき、精度が高い音声の評定が可能となる。

【0008】

30

また、本第四の発明の発音評定装置は、第一の発明の発音評定装置において、特殊音声検知部は、一の音素の評定値が所定の条件を満たすことを検知し、前記評定部は、前記特殊音声検知部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発音評定装置である。

かかる構成により、音素の置換または欠落を検出でき、精度が高い音声の評定が可能となる。

【発明の効果】

【0009】

本発明による発音評定装置によれば、無音、挿入、置換、欠落などの特殊な場合に対応した、発音の評定ができる。

40

【発明を実施するための最良の形態】

【0010】

以下、発音評定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やステップ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)

【0011】

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、かつ高速にできる発音評定装置について説明する。特に、本発音評定装置は、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

50

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。

【0012】

また、本実施の形態における発音評定装置は、例えば、語学学習や物真似練習などに利用できる。図1は、本実施の形態における発音評定装置のブロック図である。本発音評定装置は、入力受付部101、教師データ格納部102、音声受付部103、フレーム区分部104、フレーム音声データ取得部105、評定部106、出力部107を具備する。評定部106は、最適状態決定手段1061、最適状態確率値取得手段1062、評定値算出手段1063を具備する。

【0013】

入力受付部101は、発音評定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メニュー画面によるもの等、何でも良い。入力受付部101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や、メニュー画面の制御ソフトウェア等で実現され得る。

【0014】

教師データ格納部102は、教師データとして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)に基づくデータを1以上格納している。教師データは、音韻毎の隠れマルコフモデル(HMM)を連結したHMMに基づくデータであることが好適である。また、教師データは、入力される音声を構成する音素に対応するHMMを、入力順序に従って連結されているHMMに基づくデータであることが好適である。ただし、教師データは、必ずしも、音韻毎のHMMを連結したHMMに基づくデータである必要はない。教師データは、全音素のHMMの、単なる集合であっても良い。また、教師データは、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:ガウシア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。HMMに基づくデータは、例えば、フレーム毎に、状態識別子と遷移確率の情報を有する。また、HMMに基づくデータは、例えば、複数の学習対象言語を母国語として話す外国人が発声した2以上のデータから学習した(推定した)モデルでも良い。教師データ格納部102は、ハードディスクやROMなどの不揮発性の記録媒体が好適であるが、RAMなどの揮発性の記録媒体でも実現可能である。

【0015】

音声受付部103は、音声の入力を受け付ける。音声受付部103は、例えば、マイク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。

【0016】

フレーム区分部104は、音声受付部103が受け付けた音声を、フレームに区分する。フレーム区分部104は、通常、MPUやメモリ等から実現され得る。フレーム区分部104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。

【0017】

フレーム音声データ取得部105は、フレーム区分部104が区分し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部105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部105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。

【0018】

評定部106は、教師データ格納部102の教師データと、フレーム音声データ取得部

10

20

30

40

50

105が取得した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。評定部106は、通常、MPUやメモリ等から実現され得る。評定部106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。

#### 【0019】

出力部107は、評定部106の評定結果を出力する。出力部107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。出力部107は、例えば、評定部106の評定結果を視覚的に表示する。出力部107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する。出力部107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るても含まないと考えるても良い。出力部107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

#### 【0020】

評定部106を構成している最適状態決定手段1061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1061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を、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 $o_t$ 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。なお、最適状態を毛低するアルゴリズムは、例えば、Viterbial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るても良いし、連結される前のデータであり、全音韻HMMのデータと考えるても良い。

20

最適状態確率値取得手段1062は、最適状態決定手段1061が決定した最適状態における確率値を取得する。

30

#### 【0021】

評定値算出手段1063は、最適状態確率値取得手段1062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1063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1063は、例えば、最適状態確率値取得手段1062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1063は、ここでは、通常、フレームごとに評定値を算出する。

#### 【0022】

最適状態決定手段1061、最適状態確率値取得手段1062、評定値算出手段1063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1061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。

40

次に、本発音評定装置の動作について図2のフローチャートを用いて説明する。

#### 【0023】

(ステップS201)入力受付部101は、発音評定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S214に飛ぶ。

(ステップS202)音声受付部103は、音声の入力を受け付けたか否かを判断する。音声の入力を受け付ければステップS203に行き、音声の入力を受け付けなければステップS213に飛ぶ。

50

(ステップS203) フレーム区分部104は、ステップS202で受け付けた音声のデータを図示しないバッファに一時格納する。

【0024】

(ステップS204) フレーム区分部104は、バッファに一時格納した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。フレーム区分部104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部105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。

(ステップS205) フレーム音声データ取得部105は、カウンタ*i*に1を代入する。

10

【0025】

(ステップS206) フレーム音声データ取得部105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S207に行き、*i*番目のフレームが存在しなければステップS208に行く。

【0026】

(ステップS207) フレーム音声データ取得部105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ンネル数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。

20

(ステップS208) フレーム音声データ取得部105は、カウンタ*i*を1、インクリメントする。ステップS206に戻る。

【0027】

(ステップS209) 最適状態決定手段1061は、全フレームの最適状態を決定する。最適状態決定手段1061が最適状態を決定するアルゴリズムは、例えば、Viterbialgorithmによる。Viterbialgorithm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。

30

【0028】

(ステップS210) 最適状態確率値取得手段1062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1062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S211) 最適状態確率値取得手段1062は、ステップS210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。

【0029】

40

(ステップS212) 評定値算出手段1063は、ステップS211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1063が評定値を算出する関数は問わない。評定値算出手段1063は、例えば、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音声の評定値を算出する。詳細については、後述する。

【0030】

(ステップS213) 出力部107は、ステップS212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

50

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S 2 1 4で設定されるモードである。

【0031】

(ステップS 2 1 4) 音声受付部103は、タイムアウトか否かを判断する。つまり、音声受付部103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 1 4に戻る。

【0032】

(ステップS 2 1 5) 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S 2 1 6に行き、出力態様変更指示を受け付なければステップS 2 1 7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S 2 1 6) 出力部107は、ステップS 2 1 5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S 2 0 1に戻る。

(ステップS 2 1 7) 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S 2 0 1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。

以下、本実施の形態における発音評定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、発音評定装置が語学学習に利用される場合について説明する。

【0033】

まず、本発音評定装置において、図示しない手段により、ネイティブ発音の音声データベースからネイティブ発音の音韻HMMを学習しておく。ここで、音韻の種類数をLとし、1番目の音韻に対するHMMを $H_1$ とする。なお、かかる学習の処理については、公知技術であるので、詳細な説明は省略する。なお、HMMの仕様については、図3に示す。なお、HMMの仕様は、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、HMMの仕様が、他の仕様でも良いことは言うまでもない。

【0034】

そして、学習したL種類の音韻HMM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HMMを取得し、当該取得した1以上のHMM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102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。

次に、学習者が、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。

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。

次に、フレーム区分部104は、音声受付部103が受け付けた音声を、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。

【0035】

そして、フレーム音声データ取得部105は、フレーム区分部104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「 $O = o_1, o_2, \dots, o_T$ 」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ンネル数24のフィルタバンク出力を離散コサイン変換したMFC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。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

10

20

30

40

50

件を図4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。

【0036】

次に、最適状態決定手段1061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。最適状態決定手段1061が最適状態を決定するアルゴリズムは、例えば、Viterbiアルゴリズムによる。かかる場合、最適状態決定手段1061は、上記で連結したHMMを用いて最適状態を決定する。最適状態決定手段1061は、2以上のフレームの最適状態である最適状態系列を求めることとなる。

【0037】

次に、最適状態確率値取得手段1062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（ $\gamma_t(q_t^*)$ ）を算出する。なお、 $\gamma_t(q_t^*)$ は、状態 $j$ の事後確率関数 $\gamma_t(j)$ 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 $\gamma_t(j)$ は、数式2を用いて算出される。この確率値（ $\gamma_t(j)$ ）は、 $t$ 番目の特徴ベクトル $o_t$ 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】

$$DAP(t) = \gamma_t(q_t^*)$$

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（ $\gamma_t(j)$ ）は、HMMの最尤推定におけるBaum-Welchアルゴリズムの中で表れる占有度数に対応する。

【数2】

$$\begin{aligned} \gamma_t(j) &= \frac{\Pr(q_t=j, \mathbf{O} \mid \lambda^{\text{all}})}{\Pr(\mathbf{O} \mid \lambda^{\text{all}})} \\ &= \frac{\Pr(q_t=j, \mathbf{O} \mid \lambda^{\text{all}})}{\sum_{k=1}^N \Pr(q_t=k, \mathbf{O} \mid \lambda^{\text{all}})} \quad (1) \\ &= \frac{\alpha_t(j) \beta_t(j)}{\sum_{k=1}^N \alpha_t(k) \beta_t(k)} \quad (2) \end{aligned}$$

$$\alpha_t(j) = \Pr(q_t=j, \{o_1, o_2, \dots, o_t\} \mid \lambda^{\text{all}})$$

$$\beta_t(j) = \Pr(\{o_{t+1}, o_{t+2}, \dots, o_T\} \mid q_t=j, \lambda^{\text{all}})$$

数式2は、数式1を変形したものである。

【0038】

数式2において、「 $\alpha_t(j)$ 」「 $\beta_t(j)$ 」は、全部のHMMを用いて、forward-backwardアルゴリズムにより算出される。「 $\alpha_t(j)$ 」は前向き尤度、「 $\beta_t(j)$ 」は後向き尤度である。Baum-Welchアルゴリズム、forward-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。

【0039】

なお、評定部 106 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、評定部 106 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。

#### 【0040】

次に、評定値算出手段 1063 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音声の評定値を算出する。かかる場合、もし学習者の t フレーム目に対応する発声が、教師データが示す発音（例えば、正しいネイティブな発音）に近ければ、数式 2 の（2）式の分子の値が、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、結果的に最適状態の確率値（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ければ、評定値は小さくなる。なお、どのネイティブ発音にも近くないような場合は、評定値はほぼ  $1/N$  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から、通常、大きな値となり、この評定値は十分小さくなる。また、ここでは、評定値は最適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。したがって、話者性や収音環境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、学習者が正しい発音をしていれば、その変動が相殺され評定値が高いスコアを維持する。よって、評定値算出手段 1063 は、最適状態確率値取得手段 1062 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ことは、極めて好適である。

#### 【0041】

かかる評定値算出手段 1063 が算出した評定値（「DAPスコア」とも言う。）を、図 5、図 6 に示す。図 5、図 6 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 V i t e r b i アルゴリズムで求めたもの）を表わしており、図の上部に音素名を表記している。図 5 は、アメリカ人男性による英語「r i g h t」の発音の D A P 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。

#### 【0042】

図 6 は、日本人男性による英語「r i g h t」の発音の D A P 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。また、図 5 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。

#### 【0043】

そして、出力部 107 は、評定部 106 の評定結果を出力する。具体的には、例えば、出力部 107 は、図 7 に示すような態様で、評定結果を出力する。つまり、出力部 107 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、出力部 107 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 7 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。なお、本発音評定装置は、学習対象の単語（図 7 の「w o r d 1」など）や、音素（図 7 の「p 1」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。

#### 【0044】

また、図 7 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。

#### 【0045】

10

20

30

40

50

なお、図7において、発音評定装置は、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、出力部107の表示は、図8に示すような態様でも良い。図8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、出力部107の表示は、図5、図6のような表示でも良いことは言うまでもない。

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。

10

#### 【0046】

また、本実施の形態によれば、連結されたHMMである連結HMM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。

#### 【0047】

20

また、本実施の形態によれば、発音評定装置は、語学学習に利用することを主として説明したが、物真似練習などに利用できる。つまり、本発音評定装置は、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。

#### 【0048】

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8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。

#### 【0049】

30

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：ガウシア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。

#### 【0050】

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発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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、格納されているデータであり、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、前記音声受付ステップで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実行させるためのプログラム、である。

40

#### 【0051】

50

また、上記プログラムにおいて、前記教師データは、音韻毎の隠れマルコフモデル（HMM）を連結したHMMに基づくデータであり、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最適状態を決定する最適状態決定サブステップと、前記最適状態決定サブ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サブステップと、前記最適状態確率値取得サブ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サブステップを具備するプログラム、である。

【0052】

なお、上記プログラムにおいて、評定値算出サブステップにおいて、前記最適状態確率値取得サブステップで取得した確率値と、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ことは好適である。

10

（実施の形態2）

【0053】

本実施の形態における発音評定装置は、実施の形態1の発音評定装置と比較して、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、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。本実施の形態における発音評定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p-DAPと呼ぶ。

【0054】

図9は、本実施の形態における発音評定装置のブロック図である。本発音評定装置は、入力受付部101、教師データ格納部102、音声受付部103、フレーム区分部104、フレーム音声データ取得部105、評定部906、出力部107を具備する。評定部906は、最適状態決定手段1061、音韻確率値取得手段9062、評定値算出手段9063を具備する。

20

音韻確率値取得手段9062は、最適状態決定手段1061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する。ここで1以上の確率値とは、1つ以上の確率値の意味である。

【0055】

評定値算出手段9063は、音韻確率値取得手段9062が取得した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9063は、例えば、音韻確率値取得手段9062が取得した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。

30

【0056】

音韻確率値取得手段9062、および評定値算出手段9063は、通常、MPUやメモリ等から実現され得る。音韻確率値取得手段9062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、本発音評定装置の動作について図10のフローチャートを用いて説明する。図10のフローチャートにおいて、図2と異なるステップについてのみ説明する。

【0057】

40

（ステップS1001）音韻確率値取得手段9062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、音韻確率値取得手段9062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-Weichアルゴリズムの中で現れる占有度数に対応する。Baum-Weich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S1002）音韻確率値取得手段9062は、全フレームの最適状態確率値を算出する。

（ステップS1003）音韻確率値取得手段9062は、カウンタ*i*に1を代入する。

【0058】

50

(ステップS1004)音韻確率値取得手段9062は、i番目の最適状態が存在するか否かを判断する。i番目の最適状態が存在すればステップS1005に行き、i番目の最適状態が存在しなければステップS202に戻る。

(ステップS1005)音韻確率値取得手段9062は、i番目の最適状態を含む音韻全体の確率値をすべて取得する。

【0059】

(ステップS1006)評定値算出手段9063は、ステップS1005で取得した1以上の確率値に基づいて、音声の評定値を算出する。評定値算出手段9063は、例えば、音韻確率値取得手段9062が取得した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。

10

(ステップS1007)出力部107は、ステップS1006で算出した評定値を出力する。

(ステップS1008)音韻確率値取得手段9062は、カウンタiを1、インクリメントする。ステップS1004に戻る。

以下、本実施の形態における発音評定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。

【0060】

まず、学習者が、語学学習の開始の指示である動作開始指示を入力した後、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、フレーム区分部104は、音声受付部103が受け付けた音声を、短時間フレームに区分する。

20

そして、フレーム音声データ取得部105は、フレーム区分部104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「 $O = o_1, o_2, \dots, o_T$ 」を算出する。

次に、音韻確率値取得手段9062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。

【0061】

次に、最適状態決定手段1061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、各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。つまり、最適状態決定手段1061は、最適状態系列を得る。

30

【0062】

次に、音韻確率値取得手段9062は、フレーム毎に、当該フレームに対応する最適状態を含む音韻全体の確率値をすべて取得する。そして、評定値算出手段9063は、上記取得した1以上の確率値に基づいて、音声の評定値を算出する。具体的には、評定値算出手段9063は、数式3により評定値を算出する。

【数3】

$$p\text{-DAP}(t) \triangleq \sum_{j \in P(q_t^*)} \gamma_t(j)$$

40

なお、数式3において、 $P(i)$ は、i番目の状態を有しているHMMの持つ全状態の集合を示す。

【0063】

かかる評定値算出手段9063が算出した評定値(「p-DAPスコア」とも言う。)を、図11、図12に示す。図11は、アメリカ人男性による英語「right」の発音のp-DAPスコアを示す。図12は、日本人男性による英語「right」の発音のp-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。また、図11において、音素境界でスコアの落ち込みがあるものの、p-DAPは本来発音の良好なアメリカ人発音に対して、高いスコアをDAPより安定して出力して

50

いることがわかる。かかる判断は、図5のグラフと図11のグラフを比較して判断できる。また、図11において、音素 / r / のスコアが低い、この発音を聴いてみたところ / r / の発音が若干不明瞭であった。

【0064】

そして、出力部107は、算出したフレームごとの評定値を、順次出力する。かかる出力例は、図7または図8である。なお、出力部107は、図11、図12のようなグラフを出力しても良いことは言うまでもない。

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。

【0065】

また、本実施の形態によれば、連結されたHMMである連結HMM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、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。

【0066】

また、本実施の形態によれば、評定値を、各フレームにおいて、すべての音韻の中で最適な音韻の事後確率(確率値)を表しており、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、測定したい類似度を精度良く、安定して求めることができる。つまり、実施の形態1において、DAPは、全ての可能な状態に対する最適状態の事後確率を計算する。そして、総状態数Nは、通常、非常に大きくなり、入力音声によっては、評定値(DAPのスコア)が大きく低下する。つまり、例えば、あるフレームが存在する音韻内の2つの状態の過渡部にそのフレームが対応してしまえば、評定値が小さくなる。一方、音素に対する類似性を求める本実施の形態によれば、教師データの音韻との類似度(状態との類似度ではない)を測ることができ、好適である。

【0067】

さらに、本実施の形態における発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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、格納されているデータであり、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、前記音声受付ステップで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実行させるためのプログラム、である。

【0068】

また、上記プログラムにおいて、前記教師データは、音韻毎の隠れマルコフモデル(HMM)を連結したHMMに基づくデータであり、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最適状態を決定する最適状態決定サブステップと、前記最適状態決定サブステップで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する音韻確率値取得サブステップと、前記音韻確率値取得サブステップで取得した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サブステップを具備するプログラム、である。さらに、上記プログラムにおける評定値算出サブステップにおいて、前記音韻確率値取得サブステップで取得した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ことは好適である。これは、音韻確率値取得サブステップで取得した1以上の確率値の総和を音声の評定値とすることも含む。

(実施の形態3)

10

20

30

40

50

## 【 0 0 6 9 】

本実施の形態における発音評定装置は、実施の形態 1、2 の発音評定装置と比較して、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、発音区間ごとに算出される。本実施の形態における発音評定装置が算出する事後確率を、実施の形態 1 における D A P に対して t - D A P と呼ぶ。

## 【 0 0 7 0 】

図 1 3 は、本実施の形態における発音評定装置のブロック図である。本発音評定装置は、入力受付部 1 0 1、教師データ格納部 1 0 2、音声受付部 1 0 3、フレーム区分部 1 0 4、フレーム音声データ取得部 1 0 5、評定部 1 3 0 6、出力部 1 0 7 を具備する。評定部 1 3 0 6 は、最適状態決定手段 1 0 6 1、発音区間確率値取得手段 1 3 0 6 2、評定値算出手段 1 3 0 6 3 を具備する。

10

発音区間確率値取得手段 1 3 0 6 2 は、最適状態決定手段 1 0 6 1 が決定した最適状態の確率値を、発音区間毎に取得する。ここで、発音区間とは、音韻、音節、単語など、発音の一まとまりを構成する区間である。

## 【 0 0 7 1 】

評定値算出手段 1 3 0 6 3 は、発音区間確率値取得手段 1 3 0 6 2 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 1 3 0 6 3 は、例えば、発音区間確率値取得手段 1 3 0 6 2 が取得した各発音区間の 1 以上の確率値の時間平均値を、発音区間毎に算出し、1 以上の時間平均値を得て、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。

20

## 【 0 0 7 2 】

発音区間確率値取得手段 1 3 0 6 2、および評定値算出手段 1 3 0 6 3 は、通常、M P U やメモリ等から実現され得る。発音区間確率値取得手段 1 3 0 6 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 R O 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。

次に、本発音評定装置の動作について図 1 4 のフローチャートを用いて説明する。図 1 4 のフローチャートにおいて、図 1 0 と異なるステップについてのみ説明する。

（ステップ S 1 4 0 1）発音区間確率値取得手段 1 3 0 6 2 は、j に 1 を代入する。

## 【 0 0 7 3 】

（ステップ S 1 4 0 2）発音区間確率値取得手段 1 3 0 6 2 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 S 1 4 0 3 に行き、j 番目の発音区間が存在しなければステップ S 2 0 2 に行く。

30

（ステップ S 1 4 0 3）発音区間確率値取得手段 1 3 0 6 2 は、j 番目の発音区間に対応する 1 以上の最適状態の確率値をすべて取得する。

## 【 0 0 7 4 】

（ステップ S 1 4 0 4）評定値算出手段 1 3 0 6 3 は、ステップ S 1 4 0 3 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。例えば、ステップ S 1 4 0 3 で取得した 1 以上の確率値の平均値（時間平均値）を算出する。

（ステップ S 1 4 0 5）出力部 1 0 7 は、ステップ S 1 4 0 4 で算出した確率値の平均値（評定値）を出力する。

40

（ステップ S 1 4 0 6）発音区間確率値取得手段 1 3 0 6 2 は、カウンタ j を 1、インクリメントする。ステップ S 1 4 0 2 に戻る。

以下、本実施の形態における発音評定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 2 とは異なるので、その動作を中心に説明する。

## 【 0 0 7 5 】

まず、学習者が、語学学習の開始の指示である動作開始指示を入力した後、学習対象の音声を発音する。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、フレーム区分部 1 0 4 は、音声受付部 1 0 3 が受け付けた音声を、短時間フレ

50

ームに区分する。

そして、フレーム音声データ取得部 105 は、フレーム区分部 104 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「 $O = o_1, o_2, \dots, o_T$ 」を算出する。

【0076】

次に、最適状態決定手段 1061 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、各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。つまり、最適状態決定手段 1061 は、最適状態系列を得る。

次に、発音区間確率値取得手段 13062 は、各フレームの各状態の事後確率（確率値）を算出する。なお、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。

10

【0077】

そして、発音区間確率値取得手段 13062 は、発音区間に対応する 1 以上の最適状態の確率値をすべて取得する。そして、評定値算出手段 13063 は、取得した 1 以上の確率値の平均値（時間平均値）を算出する。具体的には、評定値算出手段 13063 は、数式 4 により評定値を算出する。

【数 4】

$$t\text{-DAP}(m) \triangleq \frac{\sum_{\tau \in T(q_t^*)} \gamma_{\tau}(q_{\tau}^*)}{|T(q_t^*)|}$$

20

【0078】

かかる評定値算出手段 13063 が算出した評定値（「t-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-DAP における時間平均の範囲を示す。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、出力部 107 は、算出した発音区間ごと（例えば、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。

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。

30

【0079】

また、本実施の形態によれば、連結された HMM である連結 HMM 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、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。

40

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、測定したい類似度を精度良く、安定して求めることができる。

【0080】

さらに、本実施の形態における発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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、格納されているデータであり、音韻毎の隠れマルコフモデル（HMM）に基づくデータである教師データと前記 1 以上のフ

50

フレーム音声データに基づいて、前記音声受付ステップで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実行させるためのプログラム、である。

【0081】

また、上記プログラムにおいて、前記教師データは、音韻毎の隠れマルコフモデル(HMM)を連結したHMMに基づくデータであり、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サブステップと、前記最適状態決定サブステップで決定した最適状態の確率値を、発音区間毎に取得する発音区間確率値取得サブステップと、前記発音区間確率値取得サブ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サブステップを具備するプログラム、である。

10

【0082】

さらに、評定値算出サブステップは、前記発音区間確率値取得サブステップで取得した各発音区間の1以上の確率値の時間平均値を、発音区間毎に算出し、1以上の時間平均値を得て、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ことが好適である。

(実施の形態4)

【0083】

本実施の形態における発音評定装置は、実施の形態1、2、3の発音評定装置と比較して、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発音評定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。

20

【0084】

図17は、本実施の形態における発音評定装置のブロック図である。本発音評定装置は、入力受付部101、教師データ格納部102、音声受付部103、フレーム区分部104、フレーム音声データ取得部105、評定部1706、出力部107を具備する。評定部1706は、最適状態決定手段1061、発音区間フレーム音韻確率値取得手段17062、評定値算出手段17063を具備する。

発音区間フレーム音韻確率値取得手段17062は、最適状態決定手段1061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。

30

【0085】

評定値算出手段17063は、発音区間フレーム音韻確率値取得手段17062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17063は、例えば、最適状態決定手段1061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。

【0086】

発音区間フレーム音韻確率値取得手段17062、および評定値算出手段17063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17062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。

40

次に、本発音評定装置の動作について図18のフローチャートを用いて説明する。図18のフローチャートにおいて、図14と異なるステップについてのみ説明する。

(ステップS1801) 発音区間フレーム音韻確率値取得手段17062は、カウンタkに1を代入する。

【0087】

(ステップS1802) 発音区間フレーム音韻確率値取得手段17062は、k番目の

50

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 S 1 8 0 3 に行き、k 番目のフレームが存在しなければステップ S 1 8 0 1 1 0 6 に飛ぶ。

(ステップ S 1 8 0 3) 発音区間フレーム音韻確率値取得手段 1 7 0 6 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 S 1 8 0 4) 評定値算出手段 1 7 0 6 3 は、ステップ S 1 8 0 3 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 S 1 8 0 5) 発音区間フレーム音韻確率値取得手段 1 7 0 6 2 は、k を 1、インクリメントする。ステップ S 1 8 0 2 に戻る。

#### 【0088】

(ステップ S 1 8 0 6) 評定値算出手段 1 7 0 6 3 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 1 7 0 6 3 は、例えば、最適状態決定手段 1 0 6 1 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 S 1 8 0 7) 出力部 1 0 7 は、ステップ S 1 8 0 6 で算出した評定値を出力する。

(ステップ S 1 8 0 8) 発音区間フレーム音韻確率値取得手段 1 7 0 6 2 は、j を 1、インクリメントする。ステップ S 1 4 0 2 に戻る。

以下、本実施の形態における発音評定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 3 とは異なるので、その動作を中心に説明する。

#### 【0089】

まず、学習者が、語学学習の開始の指示である動作開始指示を入力した後、学習対象の音声を発音する。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、フレーム区分部 1 0 4 は、音声受付部 1 0 3 が受け付けた音声を、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部 1 0 5 は、フレーム区分部 1 0 4 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「 $O = o_1, o_2, \dots, o_T$ 」を算出する。

次に、発音区間フレーム音韻確率値取得手段 1 7 0 6 2 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。

#### 【0090】

次に、最適状態決定手段 1 0 6 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、各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。つまり、最適状態決定手段 1 0 6 1 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)を算出と、最適状態の決定の順序は問わない。

#### 【0091】

次に、発音区間フレーム音韻確率値取得手段 1 7 0 6 2 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 1 7 0 6 3 は、各フレームの最適状態を含む音韻の全ての確率値の総和を、フレーム毎に算出する。そして、評定値算出手段 1 7 0 6 3 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 1 7 0 6 3 は、数式 5 により評定値を算出する。

10

20

30

40

【数5】

$$t\text{-}p\text{-}DAP(m) \triangleq \frac{\sum_{\tau \in T(q_t^*)} p\text{-}DAP(\tau)}{|T(q_t^*)|}$$

【0092】

かかる評定値算出手段17063が算出した評定値（「t-p-DAPスコア」とも言う。）を、図19の表に示す。図19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19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。

10

そして、出力部107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。

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。

【0093】

また、本実施の形態によれば、連結されたHMMである連結HMM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。

20

【0094】

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。

【0095】

さらに、本実施の形態における発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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、格納されているデータであり、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、前記音声受付ステップで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実行させるためのプログラム、である。

30

【0096】

また、上記プログラムにおいて、前記教師データは、音韻毎の隠れマルコフモデル（HMM）を連結したHMMに基づくデータであり、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サブステップと、前記最適状態決定サブ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サブステップと、前記発音区間フレーム音韻確率値取得サブ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サブステップを具備するプログラム、である。

40

【0097】

以上の4つの実施の形態で算出した評定値は、図20において、それぞれ（1）～（4）の類似度である。つまり、図20において、縦軸は、音韻毎の隠れマルコフモデル（HMM）を、比較対象の音素の順に連結したHMMである。図20の横軸は、入力音声の特

50

徴ベクトル系列を示す。また、図20の実線は、最適状態系列である。そして、黒丸の(1)はDAP、斜線部の(2)はp-DAP、網掛け部の(3)はt-DAPを示す。tp-DAPは、発音区間において、最適状態を含む全音韻の状態の確率値に基づいて算出される。

(実施の形態5)

【0098】

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発音評定装置について説明する。特に、本発音評定装置は、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発音評定装置である。

【0099】

また、本実施の形態における発音評定装置は、例えば、語学学習や物真似練習などに利用できる。図21は、本実施の形態における発音評定装置のブロック図である。本発音評定装置は、入力受付部101、教師データ格納部102、音声受付部103、フレーム区分部104、フレーム音声データ取得部105、特殊音声検知部2101、評定部2102、出力部108を具備する。評定部2102は、無音データ格納手段21021、無音区間検出手段21022、最適状態決定手段1061、最適状態確率値取得手段1062、評定値算出手段21023を具備する。

【0100】

特殊音声検知部2101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、特殊音声検知部2101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ことと、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。特殊音声検知部2101は、通常、MPUやメモリ等から実現され得る。特殊音声検知部2101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。

【0101】

評定部2102は、教師データ格納部102の教師データと入力音声データと特殊音声検知部2106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。評定方法の具体例は、後述する。評定部2102は、通常、MPUやメモリ等から実現され得る。評定部2102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。

【0102】

無音データ格納手段21021は、無音を示すデータであり、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21021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。

【0103】

無音区間検出手段21022は、フレーム音声データ取得部105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21021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21022は、フレーム音声データ取得部105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21022は、下記で述べる最適状態確率値取得手段1062が取得した確率値が所定の値以下であり、かつ、フレーム音声データ取得部105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。

【0104】

評定値算出手段21023は、無音区間検出手段21022が検出した無音区間を除い

10

20

30

40

50

て、かつ最適状態確率値取得手段1062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21023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21023は、例えば、最適状態確率値取得手段1062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21022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21023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21023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。

#### 【0105】

無音区間検出手段21022、評定値算出手段21023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21022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、本発音評定装置の動作について図22のフローチャートを用いて説明する。

#### 【0106】

(ステップS2201)評定値算出手段21023は、ステップS207で取得したi番目のフレーム音声データに対応する評定値(例えば、p-DAPスコア)を算出する。なお、p-DAPスコアの算出方法は、実施の形態2で述べたので、ここでの説明は省略する。

#### 【0107】

(ステップS2202)特殊音声検知部2101は、ステップS2201で算出した値が、所定の値より低いかなかを判断する。所定の値より低ければステップS2203に行き、所定の値より低くなければステップS2206に飛ぶ。

(ステップS2203)無音区間検出手段21022は、無音データと全教師データの確率値を取得する。

#### 【0108】

(ステップS2204)無音区間検出手段21022は、ステップS2203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なかを判断する。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S2205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S2206に行く。

(ステップS2205)無音区間検出手段21022は、カウンタiを1、インクリメントする。ステップS206に戻る。

(ステップS2206)出力部108は、ステップS2201で算出した評定値を出力する。

#### 【0109】

なお、図22のフローチャートにおいて、出力部108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、特殊音声を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21023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。評定部2102は、教師データと入力音声データと特殊音声検知部2101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。

#### 【0110】

また、図22のフローチャートにおいて、特殊音声検知部2101は、i番目のフレーム音声データのp-DAPスコアに基づいて特殊音声を検知したが、例えば、DAPスコアに基づいて特殊音声を検知しても良い。

#### 【0111】

以下、本実施の形態における発音評定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。

次に、学習者は、例えば、学習対象の音声を発音する。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、フレーム区分部104は、音声受付部103が受け付けた音声を、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部105は、フレーム区分部104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「 $O = o_1, o_2, \dots, o_T$ 」を算出する。

10

次に、最適状態決定手段1061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1062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。

#### 【0112】

次に、評定値算出手段21023は、例えば、最適状態決定手段1061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21023は、例えば、p-DAPスコアをフレーム毎に算出する。

20

#### 【0113】

そして、特殊音声検知部2101は、算出されたフレームに対応する評定値（p-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、特殊音声検知部2101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、特殊音声検知部2101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、特殊音声検知部2101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。

30

#### 【0114】

特殊音声検知部2101が、特殊音声を検知する場合について説明する図を図23に示す。図23(a)の縦軸は、p-DAPスコアであり、横軸はフレームを示す。図23(a)において、(V)は、V i t e r b iアライメントを示す。図23(a)において、網掛けのフレーム群におけるp-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。

#### 【0115】

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21022は、無音データ格納手段21021から無音データを取得し、当該フレーム群の各フレームのHMM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23(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（「APスコア」とも言う。）が高いことを示す。その結果、無音区間検出手段21022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23(a)において、網掛けのフレーム群におけるp-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、AP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23(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。

40

50

そして、出力部 108 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、出力部 108 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 7、図 8 である。

なお、出力部 108 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。

#### 【0116】

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。その場合、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。

10

#### 【0117】

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。

#### 【0118】

また、本実施の形態の具体例によれば、p-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（DAP、t-DAP、t-p-DAP）、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部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。

20

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。

#### 【0119】

さらに、本実施の形態における発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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入力音声データを得る入力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、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ステップにおける検知結果に基づいて、前記音声受付ステップで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実行させるためのプログラム、である。

30

#### 【0120】

また、上記プログラムは、前記特殊音声検知ステップにおいて、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サブステップと、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サブ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、前記評定ステップは、前記無音の区間を除いた区間の前記入力音声データと、前記教師データに基づいて、前記音声受付部が受け付けた音声の評定を行うことは好適である。

40

（実施の形態 6）

#### 【0121】

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発音評定装置について説明する。特に、本発音評定装置は、音韻の挿入を検知できる発音評定装置である。

#### 【0122】

また、本実施の形態における発音評定装置は、例えば、語学学習や物真似練習などに利用できる。図 24 は、本実施の形態における発音評定装置のブロック図である。本発音評

50

定装置は、入力受付部 101、教師データ格納部 102、音声受付部 103、フレーム区分部 104、フレーム音声データ取得部 105、特殊音声検知部 2401、評価部 2402、出力部 2403 を具備する。

【0123】

特殊音声検知部 2401 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価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。特殊音声検知部 2401 は、通常、MPU やメモリ等から実現され得る。特殊音声検知部 2401 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】

評価部 2402 は、特殊音声検知部 2401 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価結果を構成する。なお、評価部 2402 は、実施の形態 5 で述べたアルゴリズムにより、特殊音声検知部 2401 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ことを検知しても良い。また、評価部 2402 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、との評価結果を得ても良い。なお、実施の形態 5 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。評価部 2402 は、通常、MPU やメモリ等から実現され得る。評価部 2402 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。

【0125】

出力部 2403 は、評価部 2402 における評価結果を出力する。ここでの評価結果は、音素の挿入があった旨を示す評価結果を含む。評価結果は、音素の挿入があった旨、および評価値（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、評価値は低くなる。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。出力部 2403 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へても含まないと考へても良い。出力部 2403 は、出力デバイスのドライバソフトまたは、出力デバイスのドライバソフトと出力デバイス等で実現され得る。

次に、本発音評価装置の動作について図 25 のフローチャートを用いて説明する。

【0126】

（ステップ S2501）特殊音声検知部 2401 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 S2202 で、所定の値より低い評価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 S2507 に行き、データが格納されていなければステップ S202 に戻る。

【0127】

（ステップ S2502）特殊音声検知部 2401 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 S2507 に行き、データが格納されていなければステップ S2503 に行く。

（ステップ S2503）出力部 2403 は、ステップ S2201 で算出した評価値を出力する。

（ステップ S2504）特殊音声検知部 2401 は、カウンタ  $i$  を 1、インクリメントする。ステップ S206 に戻る。

（ステップ S2505）特殊音声検知部 2401 は、バッファに、所定の値より低い評価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 S2506）特殊音声検知部 2401 は、カウンタ  $i$  を 1、インクリメント

10

20

30

40

50

する。ステップS 2 0 6に戻る。

(ステップS 2 5 0 7) 特殊音声検知部 2 4 0 1は、カウンタ  $j$  に 1 を代入する。

【0 1 2 8】

(ステップS 2 5 0 8) 特殊音声検知部 2 4 0 1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S 2 5 0 9に行き、 $j$  番目のデータが存在しなければステップS 2 5 1 5に飛ぶ。

(ステップS 2 5 0 9) 特殊音声検知部 2 4 0 1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S 2 5 1 0) 特殊音声検知部 2 4 0 1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。

10

【0 1 2 9】

(ステップS 2 5 1 1) 特殊音声検知部 2 4 0 1は、ステップS 2 5 0 9で取得した音素とステップS 2 5 1 0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S 2 5 1 2に行き、異なる音素でなければステップS 2 5 1 4に飛ぶ。

(ステップS 2 5 1 2) 評価部 2 4 0 2は、音素の挿入があった旨を示す評価結果を構成する。

(ステップS 2 5 1 3) 特殊音声検知部 2 4 0 1は、カウンタ  $j$  を 1、インクリメントする。ステップS 2 5 0 8に戻る。

(ステップS 2 5 1 4) 出力部 2 4 0 3は、バッファ中の全データに対応する全評価値を出力する。ここで、全評価値とは、例えば、フレーム毎の  $p$ -DAPスコアである。ステップS 2 5 1 3に行く。

20

【0 1 3 0】

(ステップS 2 5 1 5) 出力部 2 4 0 3は、評価結果に「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S 2 5 1 6に行き、「挿入の旨」の情報が入っていなければステップS 2 5 1 7に行く。

(ステップS 2 5 1 6) 出力部 2 4 0 3は、評価結果を出力する。

(ステップS 2 5 1 7) 出力部 2 4 0 3は、バッファをクリアする。ステップS 2 0 6に戻る。

【0 1 3 1】

30

なお、図 2 5 のフローチャートにおいて、評価値の低いフレームが 2 つの音素に渡って存在すれば、直ちに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価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 2 5 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価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発音評価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 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。

まず、学習者が、語学学習の開始の指示である動作開始指示を入力する。

40

次に、学習者は、例えば、学習対象の音声を発音する。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、フレーム区分部 1 0 4は、音声受付部 1 0 3が受け付けた音声を、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部 1 0 5は、フレーム区分部 1 0 4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「 $O = o_1, o_2, \dots, o_T$ 」を算出する。

次に、最適状態決定手段 1 0 6 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。

50

次に、最適状態確率値取得手段1062は、数式1、2により、最適状態における確率値を算出する。

【0132】

次に、評定値算出手段21023は、例えば、最適状態決定手段1061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21023は、例えば、p-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DAPスコア等でも良い。

【0133】

そして、特殊音声検知部2101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、p-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。

10

【0134】

次に、特殊音声検知部2101は、図26に示すように、評定値(例えば、p-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25で説明した。また、図26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。

【0135】

次に、評定部2402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。そして、出力部2403は、構成した評定結果を出力する。図27は、評定結果の出力例である。なお、出力部2403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。

20

【0136】

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。その場合、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。

【0137】

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、p-DAP、t-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。

30

【0138】

さらに、本実施の形態における発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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入力音声データを得る入力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、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ステップにおける検知結果に基づいて、前記音声受付ステップで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実行させるためのプログラム、である。

40

【0139】

また、上記プログラムは、前記特殊音声検知ステップにおいて、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、前記評定ステップにおいては、前記特殊音声検知ステップにおいて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プログラム、である。

(実施の形態7)

【0140】

50

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発音評定装置について説明する。特に、本発音評定装置は、音韻の置換を検知できる発音評定装置である。

【0141】

また、本実施の形態における発音評定装置は、例えば、語学学習や物真似練習などに利用できる。図28は、本実施の形態における発音評定装置のブロック図である。本発音評定装置は、入力受付部101、教師データ格納部102、音声受付部103、フレーム区分部104、フレーム音声データ取得部105、特殊音声検知部2801、評定部2802、出力部2403を具備する。

【0142】

特殊音声検知部2801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、特殊音声検知部2801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、特殊音声検知部2801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、特殊音声検知部2801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。特殊音声検知部2801は、通常、MPUやメモリ等から実現され得る。特殊音声検知部2801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。

【0143】

評定部2802は、特殊音声検知部2801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。評定部2802は、通常、MPUやメモリ等から実現され得る。評定部2802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、本発音評定装置の動作について図29のフローチャートを用いて説明する。

【0144】

(ステップS2901)特殊音声検知部2801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S2902に行き、一の音素でなければステップS2910に行く。

【0145】

(ステップS2902)特殊音声検知部2801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t-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。

【0146】

(ステップS2903)特殊音声検知部2801は、ステップS2902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S2904に行き、所定の値より小さければステップS2910に行く。

(ステップS2904)特殊音声検知部2801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t-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。

【0147】

(ステップS2905)特殊音声検知部2801は、ステップS2904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S2906に行き、所定の値より小さければステップS2910に行く。

【0148】

(ステップS2906)特殊音声検知部2801は、予め格納されている音韻HMM(

10

20

30

40

50

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S2907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S2910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S2907) 評定部2802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S2908) 出力部2403は、ステップS2907で構成した評定結果を出力する。 10

(ステップS2909) 出力部2403は、バッファをクリアする。ステップS206に戻る。

(ステップS2910) 出力部2403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。

以下、本実施の形態における発音評定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6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。

次に、学習者は、例えば、学習対象の音声を発音する。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。 20

次に、フレーム区分部104は、音声受付部103が受け付けた音声を、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部105は、フレーム区分部104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「 $O = o_1, o_2, \dots, o_T$ 」を算出する。

次に、最適状態決定手段1061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1062は、数式1、2により、最適状態における確率値を算出する。 30

【0149】

次に、評定値算出手段21023は、例えば、最適状態決定手段1061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21023は、例えば、p-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DAPスコア等でも良い。

【0150】

そして、特殊音声検知部2101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、p-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。 40

【0151】

次に、特殊音声検知部2101は、図30に示すように、評定値(例えば、t-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、特殊音声検知部2101は、直前の音素(音素1)または直後の音素(音素2)に対する評定値(例えば、t-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して判断する。次に、特殊音声検知部2101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が一つ存在すれば、音素の置換が発生していると判断する。なお、図30において、音素 50

2において、音素の置換が発生した区間である。なお、図30において縦軸は評定値であり、当該評定値は、DAP、p-DAP、t-DAP等、問わない。

【0152】

次に、評定部2402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、出力部2403は、構成した評定結果を出力する。なお、出力部2403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。

【0153】

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。その場合、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。

10

【0154】

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、p-DAP、t-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。

【0155】

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。

20

【0156】

さらに、本実施の形態における発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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入力音声データを得る入力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、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ステップにおける検知結果に基づいて、前記音声受付ステップで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実行させるためのプログラム、である。

30

【0157】

また、上記プログラムは、前記特殊音声検知ステップにおいて、一の音素の評定値が所定の条件を満たすことを検知し、前記評定ステップにおいて、前記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または欠落があった旨を示す評定結果を構成するプログラム、である。

【0158】

また、上記プログラムでは、前記特殊音声検知ステップにおいて、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、かつ予め格納されている音韻HMMの中で、所定の値以上の評定値が得られる音韻HMMが一つ存在することを検知し、前記評定ステップにおいて、前記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。

40

（実施の形態8）

【0159】

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発音評定装置について説明する。特に、本発音評定装置は、音韻の欠落を検知できる発音評定装置である。

【0160】

また、本実施の形態における発音評定装置は、例えば、語学学習や物真似練習などに利

50

用できる。図31は、本実施の形態における発音評定装置のブロック図である。本発音評定装置は、入力受付部101、教師データ格納部102、音声受付部103、フレーム区分部104、フレーム音声データ取得部105、特殊音声検知部3101、評定部3102、出力部2403を具備する。

特殊音声検知部3101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、

特殊音声検知部3101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、特殊音声検知部3101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ことを検知しても良い。かかる場合に、特殊音声検知部3101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よ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。特殊音声検知部3101は、通常、MPUやメモリ等から実現され得る。特殊音声検知部3101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。

#### 【0161】

評定部3102は、特殊音声検知部3101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。評定部3102は、通常、MPUやメモリ等から実現され得る。評定部3102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、本発音評定装置の動作について図32のフローチャートを用いて説明する。

#### 【0162】

(ステップS3201)特殊音声検知部3101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S2911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)評定部3102は、音素の欠落があった旨を示す評定結果を構成する。

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たが、かかる条件は必須ではない。

以下、本実施の形態における発音評定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7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。

次に、学習者は、例えば、学習対象の音声を発音する。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、フレーム区分部104は、音声受付部103が受け付けた音声を、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部105は、フレーム区分部104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「 $O = o_1, o_2, \dots, o_T$ 」を算出する。

次に、最適状態決定手段1061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)

10

20

30

40

50

を決定する。

次に、最適状態確率値取得手段 1062 は、数式 1、2 により、最適状態における確率値を算出する。

【0163】

次に、評定値算出手段 21023 は、例えば、最適状態決定手段 1061 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 21023 は、例えば、p-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 DAP スコア等でも良い。

【0164】

そして、特殊音声検知部 2101 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、t-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。

【0165】

次に、特殊音声検知部 2101 は、図 33 に示すように、評定値(例えば、t-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、特殊音声検知部 2101 は、直前の音素(音素 1)または直後の音素(音素 2)に対する評定値(例えば、t-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、p-DAP、t-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。

【0166】

次に、評定部 2402 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、出力部 2403 は、構成した評定結果を出力する。なお、出力部 2403 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。

【0167】

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。その場合、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。

【0168】

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、p-DAP、t-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。

【0169】

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。

【0170】

また、本実施の形態における発音評定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、音声の入力を受け付ける音声受付ステップと、前記音声受付ステップで受け付けた音声を、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入力音声データを得る入力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、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ステップにおける検知結果に基づいて、前記音声受付ステップで受け付けた

10

20

30

40

50

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実行させるためのプログラム、である。

【0171】

また、上記プログラムは、前記特殊音声検知ステップにおいて、一の音素の評定値が所定の条件を満たすことを検知し、前記評定ステップにおいて、前記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または欠落があった旨を示す評定結果を構成するプログラム、である。

【0172】

また、上記プログラムでは、前記特殊音声検知ステップにおいて、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し、前記評定ステップにおいて、前記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が好適である。

10

【0173】

さらに、上記プログラムでは、前記特殊音声検知ステップにおいて、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することが好適である。

【0174】

また、実施の形態5から実施の形態8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。発音評定装置は、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、発音評定装置は、主として、実施の形態1から実施の形態4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良い。

20

【0175】

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。

30

【0176】

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。

【0177】

また、図34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発音評定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34は、このコンピュータシステム340の概観図であり、図35は、コンピュータシステム340のブロック図である。

40

【0178】

図34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、モニタ344と、マイク345とを含む。

【0179】

50

図17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。

10

#### 【0180】

コンピュータシステム340に、上述した実施の形態の発音評定装置の機能を実行させるプログラムは、CD-ROM3501、またはFD3502に記憶されて、CD-ROMドライブ3412またはFDドライブ3411に挿入され、さらにハードディスク3417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341に送信され、ハードディスク3417に記憶されても良い。プログラムは実行の際にRAM3416にロードされる。プログラムは、CD-ROM3501、FD3502またはネットワークから直接、ロードされても良い。

#### 【0181】

プログラムは、コンピュータ341に、上述した実施の形態の発音評定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。

20

#### 【0182】

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。

なお、上記プログラムにおいて、ハードウェアによって行われる処理、例えば、出力ステップにおけるディスプレイなどで行われる処理(ハードウェアでしか行われない処理)は含まれない。

30

また、上記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】

#### 【0183】

以上のように、本発明にかかる発音評定装置は、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という効果を有し、語学学習装置や物真似練習装置等として有用である。

40

#### 【図面の簡単な説明】

#### 【0184】

【図1】実施の形態1における発音評定装置のブロック図

【図2】同発音評定装置の動作について説明するフローチャート

【図3】同HMMの仕様を説明する図

【図4】同音声分析条件を説明する図

【図5】同評定値算出手段が算出した評定値を示すグラフを示す図

【図6】同評定値算出手段が算出した評定値を示すグラフを示す図

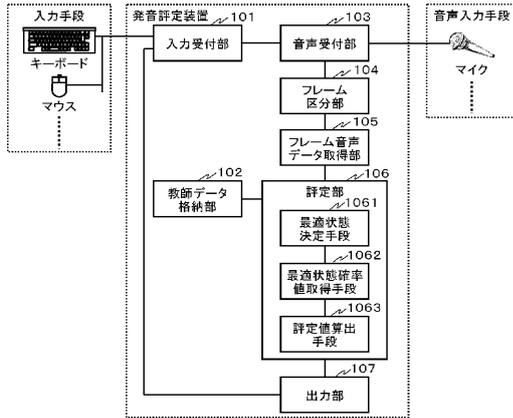
【図7】同出力部が出力する出力態様を示す図

【図8】同出力部が出力する出力態様を示す図

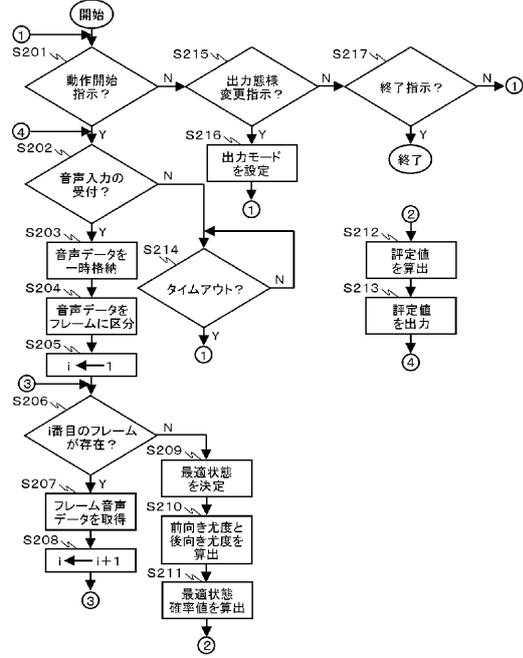
50

【図 9】実施の形態 2 における発音評定装置のブロック図	
【図 10】同発音評定装置の動作について説明するフローチャート	
【図 11】同評定値算出手段が算出した評定値を示すグラフを示す図	
【図 12】同評定値算出手段が算出した評定値を示すグラフを示す図	
【図 13】実施の形態 3 における発音評定装置のブロック図	
【図 14】同発音評定装置の動作について説明するフローチャート	
【図 15】同評定値算出手段が算出した評定値を示す図	
【図 16】同出力部が出力する出力態様を示す図	
【図 17】実施の形態 4 における発音評定装置のブロック図	
【図 18】同発音評定装置の動作について説明するフローチャート	10
【図 19】同評定値算出手段が算出した評定値を示す図	
【図 20】4 つの実施の形態で算出した評定値について説明する図	
【図 21】実施の形態 5 における発音評定装置のブロック図	
【図 22】同発音評定装置の動作について説明するフローチャート	
【図 23】同特殊音声の検知について説明する図	
【図 24】実施の形態 6 における発音評定装置のブロック図	
【図 25】同発音評定装置の動作について説明するフローチャート	
【図 26】同特殊音声の検知について説明する図	
【図 27】同評定結果の出力例を示す図	
【図 28】実施の形態 7 における発音評定装置のブロック図	20
【図 29】同発音評定装置の動作について説明するフローチャート	
【図 30】同特殊音声の検知について説明する図	
【図 31】実施の形態 8 における発音評定装置のブロック図	
【図 32】同発音評定装置の動作について説明するフローチャート	
【図 33】同特殊音声の検知について説明する図	
【図 34】同発音評定装置を構成するコンピュータシステムの概観図	
【図 35】同発音評定装置を構成するコンピュータのブロック図	
【符号の説明】	
【0185】	
101 入力受付部	30
102 教師データ格納部	
103 音声受付部	
104 フレーム区分部	
105 フレーム音声データ取得部	
106、906、1306、1706、2102、2402、2802、3102 評 定部	
107、2403 出力部	
1061 最適状態決定手段	
1062 最適状態確率値取得手段	
1063、9063、13063、17063、21023 評定値算出手段	40
9062 音韻確率値取得手段	
13062 発音区間確率値取得手段	
17062 発音区間フレーム音韻確率値取得手段	
2101、2401、2801、3101 特殊音声検知部	
21021 無音データ格納手段	
21022 無音区間検出手段	

【図1】



【図2】



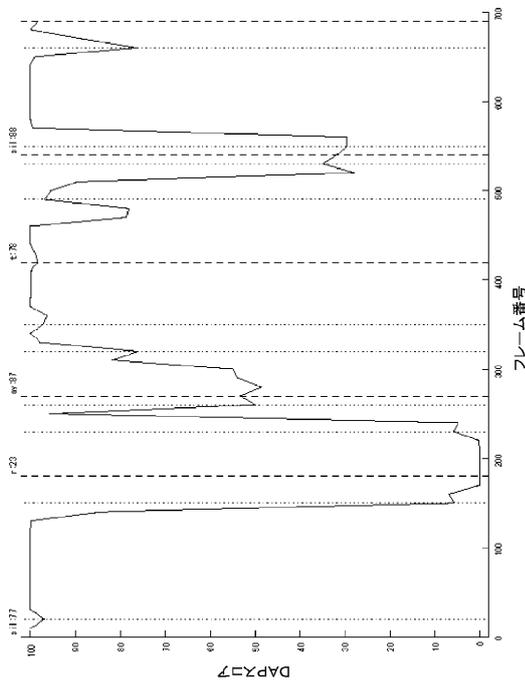
【図3】

HMMタイプ	混合ガウス型 mono phone (対角共分散行列)
状態数(モデル毎)	3
混合数(状態毎)	8
HMM総数	44
総状態数	132

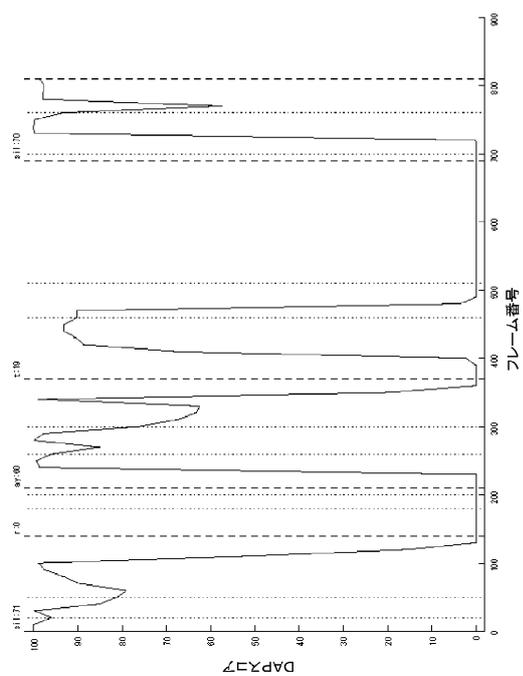
【図4】

サンプリング周波数	22.05 KHz
プリエンファシス	1-0.97z <sup>-1</sup>
窓関数	Hamming窓
分析フレーム長	25 msec
フレーム周期	10 msec
特徴パラメータ	MFCC (39次元)

【図5】

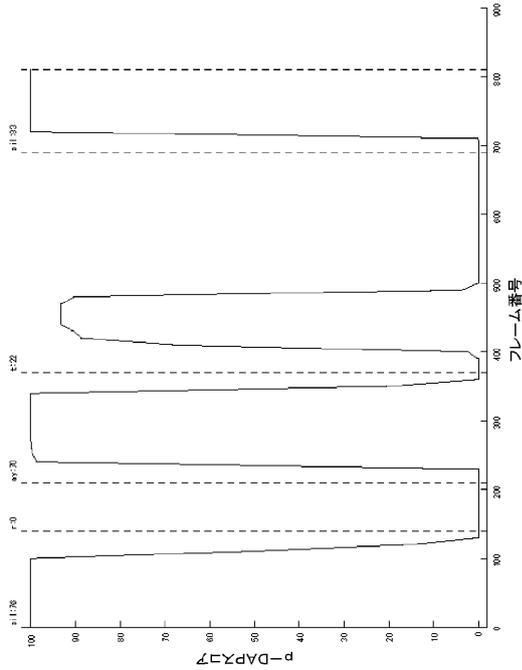


【図6】

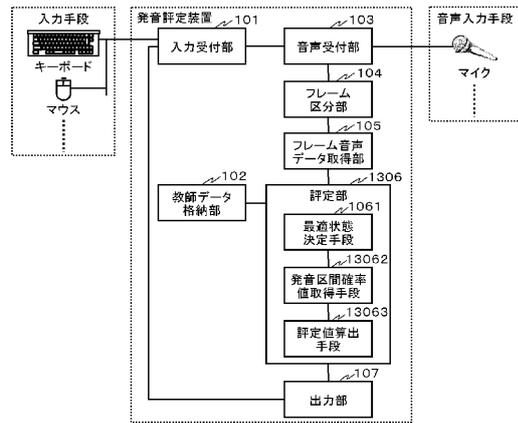




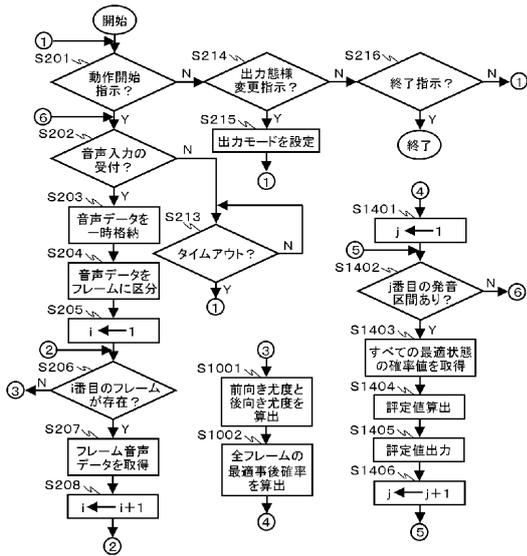
【図12】



【図13】



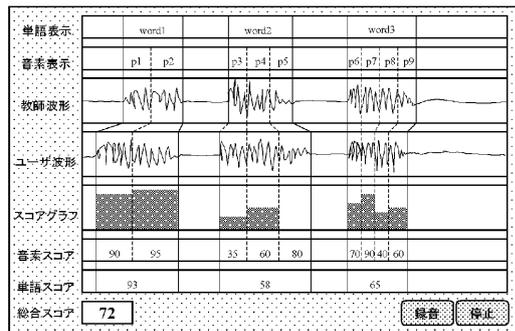
【図14】



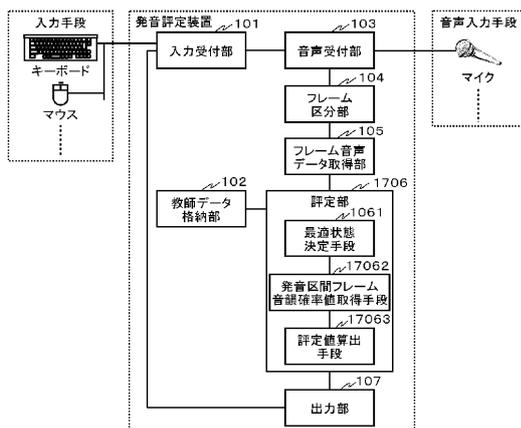
【図15】

	Phoneme			Word
	/r/	/ay/	/t/	/right/
アメリカ人男性	23	87	78	68
日本人男性	0	60	19	29

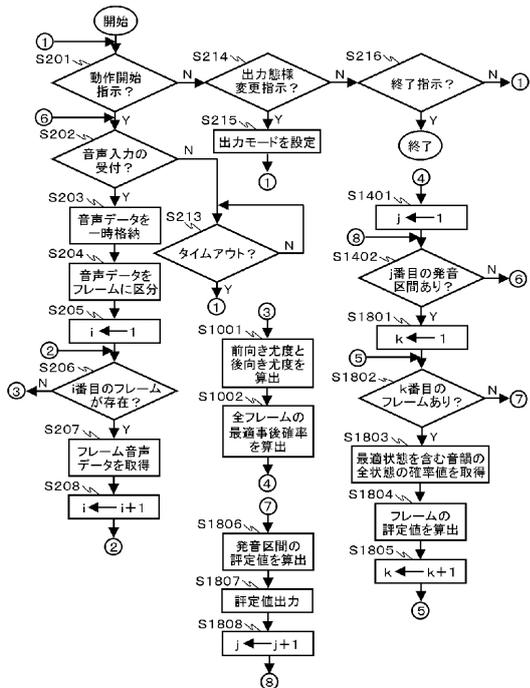
【図16】



【図17】



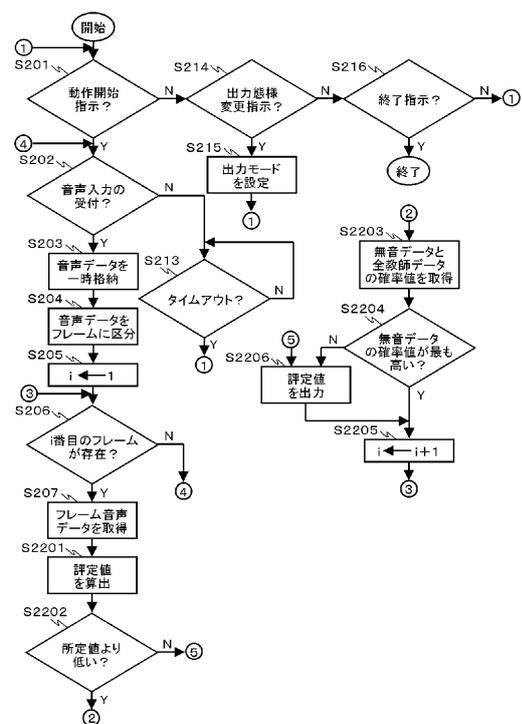
【図18】



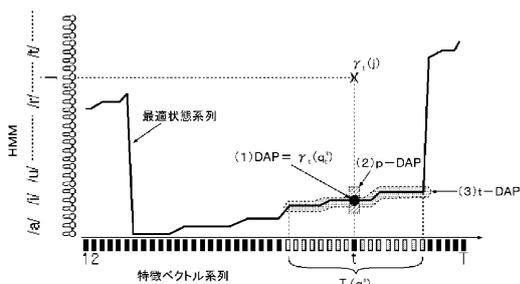
【図19】

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

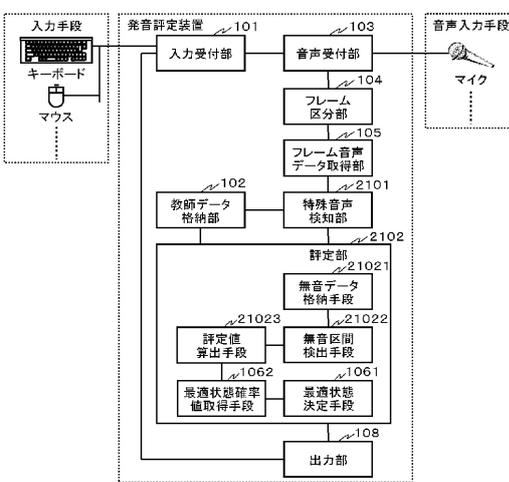
【図22】



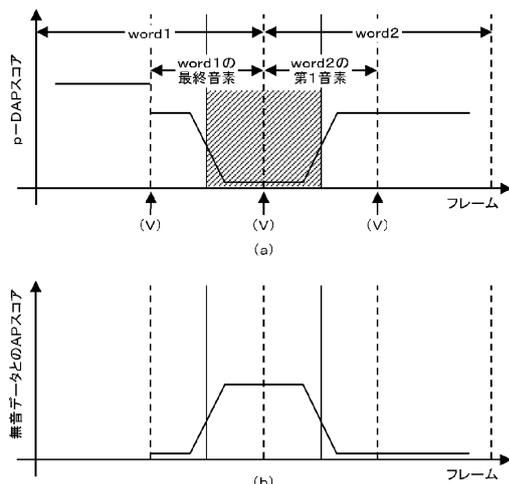
【図20】



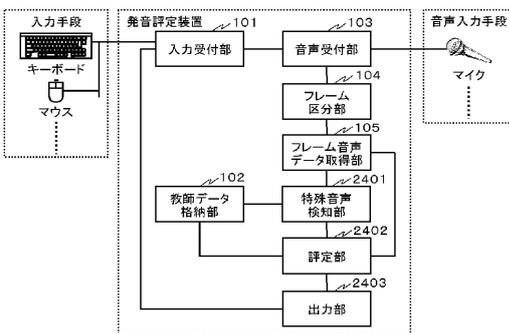
【図21】



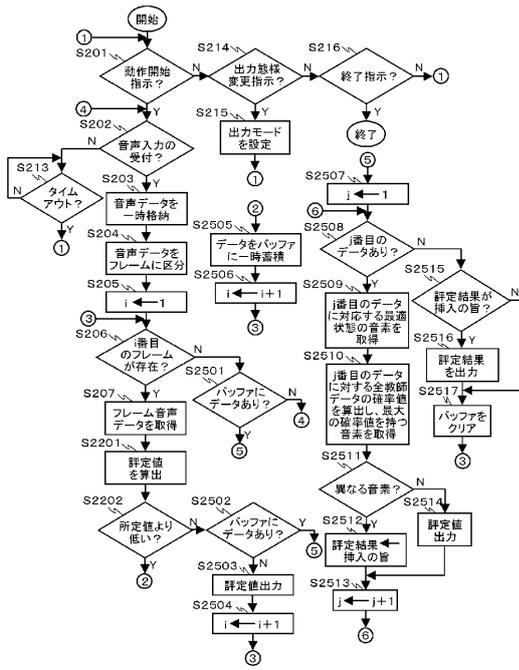
【図23】



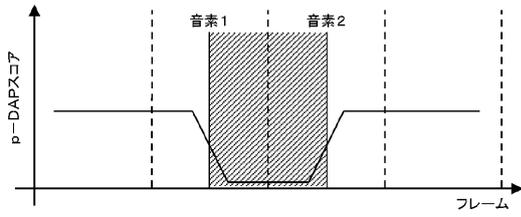
【図24】



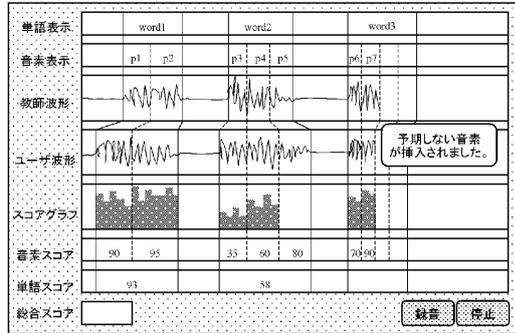
【図25】



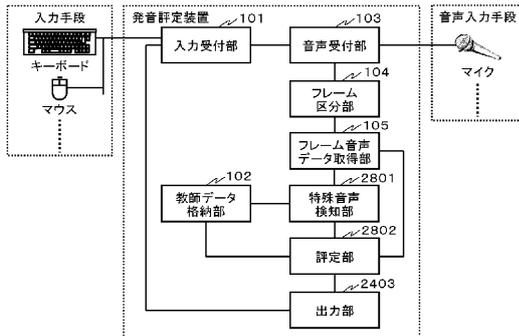
【図26】



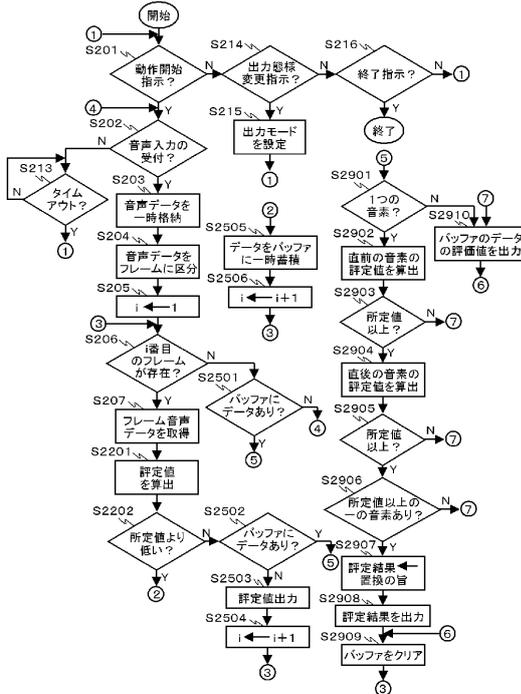
【図27】



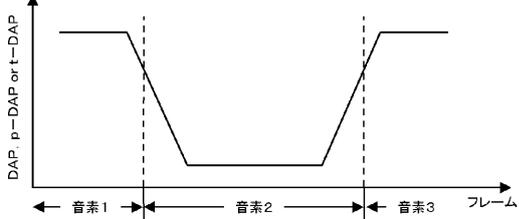
【図28】



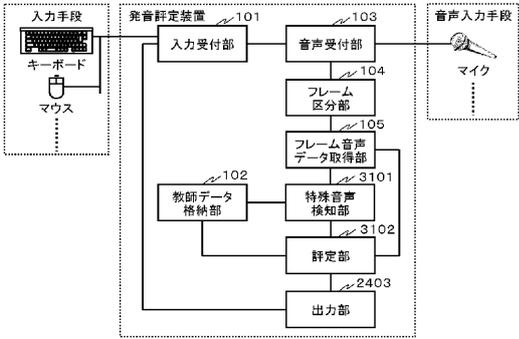
【図29】



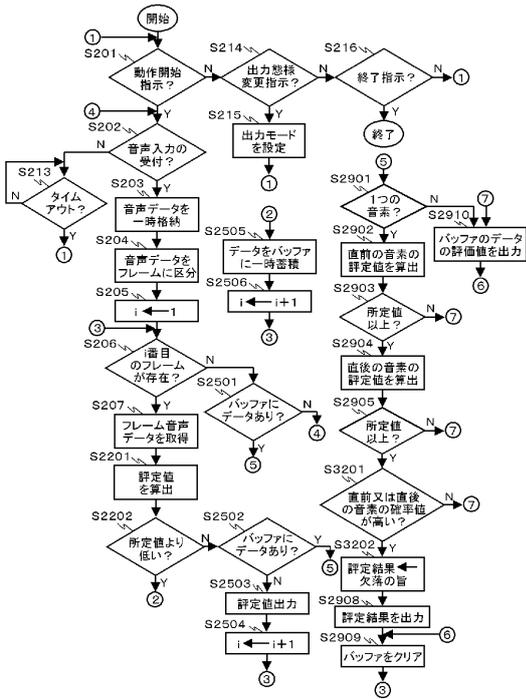
【図 30】



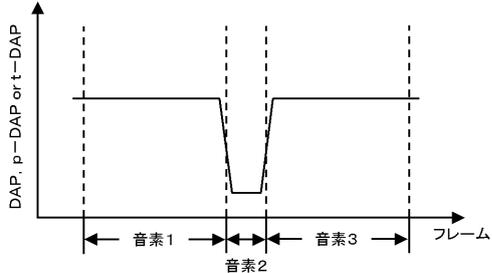
【図 31】



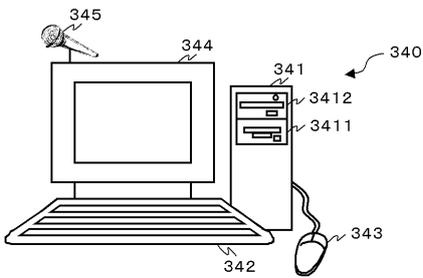
【図 32】



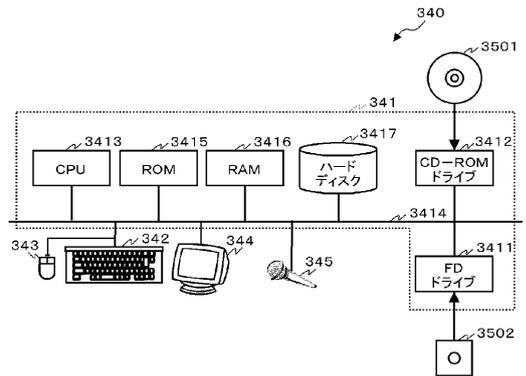
【図 33】



【図 34】



【図 35】



---

フロントページの続き

審査官 植田 泰輝

- (56)参考文献 特開2001-265211(JP,A)  
特開平9-258765(JP,A)  
特開平6-110494(JP,A)  
特開2003-345380(JP,A)  
特表平7-503559(JP,A)  
特開2004-117530(JP,A)

(58)調査した分野(Int.Cl., DB名)

G09B 5/00 - 7/12, 19/00 - 19/26  
G10L 15/00 - 15/28