

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4202124号
(P4202124)

(45) 発行日 平成20年12月24日(2008.12.24)

(24) 登録日 平成20年10月17日(2008.10.17)

(51) Int.Cl. F I
G 1 0 L 15/06 (2006.01) G 1 0 L 15/06 5 0 0 N
 G 1 0 L 15/06 3 0 0 C

請求項の数 26 (全 16 頁)

(21) 出願番号	特願2002-512966 (P2002-512966)	(73) 特許権者	595020643
(86) (22) 出願日	平成13年7月11日(2001.7.11)		クォアルコム・インコーポレイテッド
(65) 公表番号	特表2004-504641 (P2004-504641A)		QUALCOMM INCORPORATED
(43) 公表日	平成16年2月12日(2004.2.12)		ED
(86) 国際出願番号	PCT/US2001/022009		アメリカ合衆国、カリフォルニア州 92
(87) 国際公開番号	W02002/007145		121-1714、サン・ディエゴ、モア
(87) 国際公開日	平成14年1月24日(2002.1.24)		ハウス・ドライブ 5775
審査請求日	平成16年8月18日(2004.8.18)	(74) 代理人	100058479
(31) 優先権主張番号	09/615,572		弁理士 鈴江 武彦
(32) 優先日	平成12年7月13日(2000.7.13)	(74) 代理人	100091351
(33) 優先権主張国	米国 (US)		弁理士 河野 哲
		(74) 代理人	100088683
			弁理士 中村 誠
		(74) 代理人	100109830
			弁理士 福原 淑弘

最終頁に続く

(54) 【発明の名称】 話者独立音声認識システムのための音声テンプレートを構成するための方法及び装置

(57) 【特許請求の範囲】

【請求項 1】

話者独立音声認識システムにおける使用のための音声テンプレートを生成する方法であって、前記方法は、

第1の複数の発話の各発話を細分化して、各発話に対して複数の時間クラスタセグメントを生成し、各時間クラスタセグメントはスペクトル手段によって表され、

前記第1の複数の発話のすべてに対して前記複数のスペクトル手段を量子化して複数のテンプレートベクトルを生成し、

動的時間ゆがみ演算を使用して前記複数のテンプレートベクトルの各1つと第2の複数の発話とを比較して、少なくとも1つの比較結果を生成し、

前記少なくとも1つの比較結果が少なくとも1つの所定のしきい値を越えるならば、前記第1の複数の発話と前記複数のテンプレートベクトルとのマッチングを行って、最適なマッチングパス結果を生成し、

前記最適なマッチングパス結果に従って時間上で前記第1の複数の発話を分割し、

前記少なくとも1つの比較結果が少なくとも1つの所定のしきい値を越えないようになるまで前記量子化、前記比較、前記マッチング、前記分割を反復することを具備する方法。

【請求項 2】

前記比較は分散測定値を計算することを含む請求項1に記載の方法。

【請求項 3】

前記比較は精度測定値を計算することを含む請求項 1 に記載の方法。

【請求項 4】

前記比較はまず分散測定値を計算し、次に前記分散測定値が第 1 の所定のしきい値を越えないならば、精度測定値を計算することを含む請求項 1 に記載の方法。

【請求項 5】

前記分散測定値が前記第 1 の所定のしきい値を越えるかあるいは前記精度測定値が第 2 の所定のしきい値を越えるならば、前記マッチングは前記第 1 の発話と前記複数のテンプレートベクトルとのマッチングをとることを含む請求の範囲第 4 項に記載の方法。

【請求項 6】

前記マッチングは、動的時間ひずみの計算を実行することを含む請求項 1 に記載の方法

10

【請求項 7】

前記マッチング及び前記分割は、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えたときに、K - 平均細分化の手法を用いて実行されることを含む請求項 1 に記載の方法。

【請求項 8】

前記第 1 の発話の終了点を検出することをさらに含む請求項 1 に記載の方法。

【請求項 9】

話者独立音声認識システムにおける使用のための音声テンプレートを生成するように構成された装置であって、前記装置は、

20

第 1 の複数の発話の各発話を細分化して、各発話に対する複数の時間クラスタセグメントを生成する手段であって、各時間クラスタセグメントはスペクトル手段によって表される手段と、

前記第 1 の複数の発話のすべてに対して前記複数のスペクトル手段を量子化して、複数のテンプレートベクトルを生成する手段と、

前記複数のテンプレートベクトルの各 1 つと第 2 の複数の発話とを比較して、少なくとも 1 つの比較結果を生成するために動的時間ゆがみ演算を使用する手段と、

前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えるならば、前記第 1 の複数の発話と前記複数のテンプレートベクトルとのマッチングを行って、最適なマッチングパス結果を生成する手段と、

30

前記最適なマッチングパス結果に従って時間上で前記第 1 の複数の発話を分割するための手段と、

前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えなくなるまで前記量子化、前記比較、前記マッチング及び前記分割を反復する手段と、

を具備する装置。

【請求項 10】

話者独立音声認識システムにおける使用のための音声テンプレートを生成するように構成された装置であって、前記装置は、

第 1 の複数の発話の各発話を細分化して、各発話に対する複数の時間クラスタセグメントを生成するように構成される細分化ロジックであって、各時間クラスタセグメントはスペクトル手段によって表される細分化ロジックと、

40

前記細分化ロジックに結合され、前記第 1 の複数の発話のすべてに対する前記複数のスペクトル手段を量子化して、複数のテンプレートベクトルを生成するように構成された量子化器と、

前記量子化器に結合され、動的時間ゆがみ演算を使用して前記複数のテンプレートベクトルの各 1 つと、第 2 の複数の発話とを比較して、少なくとも 1 つの比較結果を生成するように構成されたコンパジェンス試験器と、

前記量子化器とコンパジェンス試験器とに結合され、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えるならば、前記第 1 の複数の発話と前記複数のテンプレートベクトルとのマッチングをとって最適なマッチングパス結果を生成し、前記

50

最適なマッチングパス結果に従って時間上で前記第 1 の複数の発話を分割するように構成された分割ロジックと

を具備し、

前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えなくなるまで前記量子化、前記比較、前記マッチング、そして前記分割を反復するように構成されている装置。

【請求項 1 1】

前記少なくとも 1 つの比較結果は分散測定値である請求項 1 0 に記載の装置。

【請求項 1 2】

前記少なくとも 1 つの比較結果は精度測定値である請求項 1 0 に記載の装置。

10

【請求項 1 3】

前記少なくとも 1 つの比較結果は分散測定値と精度測定値であり、前記コンバージョン試験器はまず前記分散測定値を計算し、次に、この分散測定値が第 1 の所定のしきい値を越えないならば、前記精度測定値を計算するように構成された請求項 1 0 に記載の装置。

【請求項 1 4】

前記マッチングは、前記分散測定値が前記第 1 の所定のしきい値を越えるかあるいは前記精度測定値が第 2 の所定のしきい値を越えるならば、前記第 1 の発話と前記複数のテンプレートベクトルとのマッチングをとることを含む請求項 1 3 に記載の装置。

【請求項 1 5】

前記分割ロジックは動的時間ゆがみ計算を実行するように構成される請求項 1 0 に記載の装置。

20

【請求項 1 6】

前記分割ロジックは、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えたときに実行されるように構成された K - 平均細分化ロジックを含む請求項 1 0 に記載の装置。

【請求項 1 7】

前記細分化ロジックに結合され、前記第 1 の発話の終了点を検出するように構成された終了点検出器をさらに含む請求項 1 0 に記載の装置。

【請求項 1 8】

話者独立音声認識システムにおける使用のための音声テンプレートを生成するように構成された装置であって、前記装置は、

30

プロセッサと、

前記プロセッサに結合され、前記プロセッサによって実行可能な一組の命令を含み、各発話に対する複数の時間クラスタセグメントを生成するために第 1 の複数の発話の各発話を分割し、各時間クラスタセグメントは手段によって表され、複数のテンプレートベクトルを生成するために前記第 1 の複数の発話のすべてに対する前記複数のスペクトル手段を量子化し、少なくとも 1 つの比較結果を生成するために動的時間ゆがみ演算を使用して前記複数のテンプレートベクトルの各 1 つと第 2 の複数の発話とを比較し、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えるならば、最適なマッチングパス結果を生成するために、前記第 1 の複数の発話と前記複数のテンプレートベクトルとのマッチングをとり、前記最適なマッチングパス結果に従って時間において前記第 1 の複数の発話を分割し、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えなくなるまで前記量子化、前記比較、前記マッチングそして前記分割を反復する記憶媒体と、

40

を具備する装置。

【請求項 1 9】

前記少なくとも 1 つの比較結果は分散測定値である請求項 1 8 に記載の装置。

【請求項 2 0】

前記少なくとも 1 つの比較結果は精度測定値である請求項 1 8 に記載の装置。

50

【請求項 2 1】

前記少なくとも 1 つの比較結果は、分散測定値と精度測定値であり、前記一組の命令は、前記分散測定値をまず計算し、次に、前記分散測定値が第 1 の所定のしきい値を越えないならば、前記精度測定値を計算するために前記プロセッサによって実行可能である請求項 1 8 に記載の装置。

【請求項 2 2】

前記一組の命令はさらに、前記分散測定値が前記第 1 の所定のしきい値を越えるかあるいは前記精度測定値が第 2 の所定のしきい値を越えるならば、前記第 1 の発話と前記複数のテンプレートベクトルとのマッチングをとるために前記プロセッサによって実行可能である請求項 2 1 に記載の装置。

10

【請求項 2 3】

前記一組の命令は、動的時間ゆがみ計算を実行することによって、前記第 1 の発話に一致するように構成された分割ロジックと前記複数のテンプレートベクトルとのマッチングをとるために、前記プロセッサによって実行可能である請求項 1 8 に記載の装置。

【請求項 2 4】

前記一組の命令は、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えたときに K - 平均音声細分化計算を実行することによって、前記第 1 の発話を分割するために、前記プロセッサによって実行可能である請求項 1 8 に記載の装置。

【請求項 2 5】

前記一組の命令はさらに、前記第 1 の発話の終了点を検出するために前記プロセッサによって実行可能である請求項 1 8 に記載の装置。

20

【請求項 2 6】

各発話に対する複数の時間クラスタセグメントを生成するために第 1 の複数の発話の各発話を分割し、各時間クラスタセグメントはスペクトル手段によって表され、

複数のテンプレートベクトルを生成するために前記第 1 の複数の発話のすべてに対して前記複数のスペクトル手段を量子化し、

少なくとも 1 つの比較結果を生成するために動的時間ゆがみ演算を使用して前記複数のテンプレートベクトルの各 1 つと第 2 の複数の発話とを比較し、

前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えるならば、前記第 1 の複数の発話と前記複数のテンプレートベクトルとのマッチングをとって最適なマッチングパス結果を生成し、

30

前記最適なマッチングパス結果に従って時間上で前記第 1 の複数の発話を分割し、

前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えなくなるまで、前記量子化、前記比較、前記マッチングそして前記分割を反復するために、

プロセッサによって実行可能な一組の命令を含むプロセッサ読取可能な媒体。

【発明の詳細な説明】

【0001】

発明の背景

1. 発明の分野

本発明は概して通信の分野に関し、より詳細には、話者独立音声認識システムのための音声テンプレートに関する。

40

【0002】

2. 背景

音声認識 (V R) は、ユーザまたはユーザ発声コマンドを認識して機械とのヒューマンインタフェースを確立するための、シミュレートされた知能を当該機械に与える最も重要な技術の 1 つを代表するものである。音響音声信号から言語メッセージを再生するための技術を使用するシステムは音声認識装置と呼ばれる。“音声認識装置”の術語は、ここでは、概して話者インタフェース駆動による装置を意味するのに使用される。音声認識装置は概して、音響プロセッサとワードデコーダとを備える。音響プロセッサは、到来する生の音声の V R を行うのに必要な、一連の情報保持特徴またはベクトルを抽出する。ワード

50

デコーダは、入力音声に対応する一連の言語ワードなどの意味のある所望の出力フォーマットを得るために、前記一連の特徴あるいはベクトルを復号する。

【 0 0 0 3 】

音響プロセッサは、音声認識装置内のフロントエンド音声解析サブシステムを代表する。入力音声信号に回答して、音響プロセッサは、時変音声信号を特徴付けるための適切な表示を提供する。音響プロセッサは、背景ノイズ、チャンネルひずみ、話者特性、そして話しかたなどの不要な情報を捨てる。効率の良い音響処理は、音声認識に強化された音響識別力を与える。この点において、解析するべき有益な特性は、短時間スペクトル包絡線である。短時間スペクトル包絡線の特徴付けるための2つの通常使用されるスペクトル解析技術は、線形予測符号化(LPC)及びフィルタバンクに基づくスペクトルモデル化である。一般的なLPC技術は、本発明の譲受人に譲渡され、ここに参照として組み込まれている米国特許第5414796号及び、ここにその全体が参照として組み込まれているL.B. Rabiner&R.W.Schafer、音声信号のデジタル処理、396-453(1978)に開示されている。

10

【 0 0 0 4 】

VR(一般的に音声認識とも呼ばれる)の使用が安全性の観点から重要になってきている。例えば、VRは、ワイヤレス電話キーパッド上の釦を押圧するというマニュアル作業に取って代わるのに使用される。このことは、自動車の運転中に電話をかけたいときに特に重要である。VRなしに電話を使用する場合、運転者は一方の手をハンドルから離して電話呼をダイヤルするために釦を押している間電話キーパッドを見なければならない。これらの行為は、自動車事故の可能性を増大させる。音声駆動による電話(すなわち、音声認識用に設計された電話)は、運転者が道路を連続的に眺めながら電話をかけることを可能にする。さらに、ハンドフリーなカーキットシステムは、運転者が電話開始の間両手をハンドルに維持することを可能にする。

20

【 0 0 0 5 】

音声認識装置は、話者依存の装置と話者独立の装置とに区別される。話者依存装置はより知られているが、特定のユーザからのコマンドを認識するように学習する。これとは対照的に、話者独立装置は任意のユーザからの音声コマンドを受け入れることができる。所定のVRシステムのパフォーマンスを増大するために、話者依存であるか話者独立であるかにかかわらず、システムに有効なパラメータを備えさせるために学習が必要である。言い換えると、システムはそれが最適に機能できるようになる前に学習する必要がある。

30

【 0 0 0 6 】

話者依存VR装置は概して、トレーニング段階と認識段階の2つの段階において動作する。トレーニング段階において、VRシステムはシステムのボキャブラリ内のワードのそれぞれを一度あるいは二度(通常は二度)話すことをユーザに促し、これにより、システムは特定のワードまたはフレーズに対するユーザ音声の特徴を学ぶことができる。ハンドフリーカーキットに対する例示的なボキャブラリは、キーパッド上の数字、“発呼”、“送信”、“ダイヤル”、“キャンセル”、“クリア”、“追加”、“削除”、“ヒストリ”、“プログラム”、“イエス”、“ノー”のキーワード、そして、通常、仕事仲間、友達あるいは家族のメンバ、と呼ばれる特定の数の人の名前を含む。学習がいったん完了したならば、ユーザは学習されたキーワードを話すことによって認識段階において発呼を開始することができ、VR装置は話された音声を(テンプレートとして記憶されている)予め学習された音声と比較し、最善のマッチングを行うことによって認識する。例えば、“ジョン”という名前が学習された名前のものであるならば、ユーザは“ジョンを呼べ”というフレーズを発音することによってジョンに対する発呼を開始する。VRシステムは“呼べ”及び“ジョン”というワードを認識し、ユーザが予めジョンの電話番号として入力した番号をダイヤルする。

40

【 0 0 0 7 】

話者独立VR装置は、所定のサイズの予め記録されたボキャブラリ(例えば、ある種の制御ワード、0から9までの数、イエス及びノー)を含む学習テンプレートを使用する。多数の話者(例えば100)がボキャブラリ内の各ワードを発音して記録しなければならな

50

い。

【 0 0 0 8 】

従来、話者独立 V R テンプレートは、第 1 の組の話者（概して 1 0 0 の話者）により話されたワードを含む試験データベースと、第 2 の組の話者（第 1 の組と同じ程度の話者）により話された同じワードを含む学習データベースとを比較することによって構成される。一人のユーザにより話された 1 つのワードは概して発話(utterance)と呼ばれる。トレーニングデータベースの各発話はまず正規化され、試験データベースの発話によりコンバージェンス (convergence) が試験される前に、量子化 (概して既知の技術に従ったベクトル量子化) される。しかしながら、時間正規化技術は、以前のフレームとは最大の相異をもつ個々のフレーム (発話の周期的セグメント) からのみ獲得される情報に依存する。所定の発話における情報の多くを使用する話者独立 V R テンプレートを構築するための方法を提供することが好ましい。発話のタイプに基づいて話者独立 V R テンプレートを構築するための従来の技術の精度またはコンバージェンスをさらに増大させることが望ましい。すなわち、増大した精度を提供するとともに発話における大量の情報を使用する話者独立音声認識テンプレートを構築する方法に対する要求が存在する。

10

【 0 0 0 9 】

発明の要約

本発明は、増大した精度を提供するとともに発話における大量の情報を使用する話者独立音声認識テンプレートを構築する方法に関する。すなわち、本発明の一側面において、話者独立音声認識システムにおける使用のための音声テンプレートを生成する方法が提供される。この方法は好ましくは、第 1 の複数の発話の各発話を分割して、各発話に対する複数の時間クラスタセグメントを生成し、各時間クラスタセグメントはスペクトル手段により表され、前記第 1 の複数の発話のすべてに対する前記複数のスペクトル手段を量子化して複数のテンプレートベクトルを生成し、前記複数のテンプレートベクトルの各々と第 2 の複数の発話とを比較して少なくとも 1 つの比較結果を生成し、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値を越えるならば、前記第 1 の複数の発話と前記複数のテンプレートベクトルとのマッチングを行って、最適なマッチングパス結果を生成し、前記最適なマッチングパス結果に従って、時間上で前記第 1 の複数の発話を分割し、前記少なくとも 1 つの比較結果が少なくとも 1 つの所定のしきい値のどれをも越えなくなるまで前記量子化、前記比較、前記マッチング、前記分割を反復する。

20

30

【 0 0 1 0 】

好ましい実施形態の詳細な説明

一実施形態によれば、図 1 に示すように、話者独立音声認識のための音声テンプレートを構築して実行するためのシステム 1 0 は、話者独立テンプレート構築サブシステム 1 2 及び音声認識サブシステム 1 4 を含む。話者独立テンプレート構築サブシステム 1 2 は、音声認識サブシステム 1 4 に結合される。

【 0 0 1 1 】

話者独立音声テンプレートは、図 4 ~ 図 6 を参照して後述するように、話者独立テンプレート構築サブシステム 1 2 により構築される。テンプレートは、図 2 - 3 を参照して後述するように、ユーザからの入力音声を認識するにおいて使用される音声認識サブシステム 1 4 に供給される。

40

【 0 0 1 2 】

一実施形態によれば、図 2 に示すように、音声認識サブシステム 1 0 0 は、アナログ/デジタル変換器 (A/D) 1 0 2、フロントエンド音響プロセッサ 1 0 4、特徴抽出器 1 0 6、音声テンプレートデータベース 1 0 8、パターン比較ロジック 1 1 0、決定ロジック 1 1 2 を含む。特定の実施形態において、音響プロセッサ 1 0 4 及び特徴抽出器 1 0 6 は 1 つのデバイス例えばパラメータ抽出器として実現される。一実施形態において音響プロセッサ 1 0 4 は周波数解析モジュール 1 1 4 を含む。一実施形態において特徴抽出器 1 0 6 は、終了点検出器 1 1 6、時間クラスタリング音声細分化モジュール 1 1 8、音声レベル正規化器 1 2 0 を含む。

50

【 0 0 1 3 】

A / D 1 0 2 は、音響プロセッサ 1 0 4 に結合されている。音響プロセッサ 1 0 4 は、特徴抽出器 1 0 6 に結合される。一実施形態において、特徴抽出器 1 0 6 内で、終了点検出器 1 1 6 は、振幅量子化器 1 2 0 に結合された時間クラスタリング音声細分化モジュール 1 1 8 に結合されている。特徴抽出器 1 0 6 はパターン比較ロジック 1 1 0 に結合されている。パターン比較ロジック 1 1 0 はテンプレートデータベース 1 0 8 及び決定ロジック 1 1 2 に結合されている。

【 0 0 1 4 】

音声認識サブシステム 1 0 0 は例えばワイヤレス電話またはハンドフリーカーキット内に配置されている。(図示せぬ)ユーザがワードまたはフレーズを発音すると音声信号を生成する。音声信号は、従来の変換器(図示せぬ)によって電気音声信号 $s(t)$ に変換される。音声信号 $s(t)$ は A / D 1 0 2 に供給されて当該音声信号が例えばパルス符号変調(PCM)、A-law、あるいは μ -law などの既知のサンプリング方法に従ってデジタル音声サンプル $s(n)$ に変換される。

【 0 0 1 5 】

音声サンプル $s(n)$ はパラメータ決定のために音響プロセッサ 1 0 4 に供給される。音響プロセッサ 1 0 4 は入力音声信号 $s(t)$ の特性をまねた一連のパラメータを生成する。パラメータは、例えば音声符号器符号化、離散フーリエ変換(DFT)に基づくケプストラム係数(例えば高速フーリエ変換(FFT)に基づくケプストラム係数)、線形予測係数(LPC)あるいはパークスケール解析、を含む既知の音声パラメータ決定技術に従って決定される。このような方法は前記した米国特許第 5 4 1 4 7 9 6 及び Lawrence Rabiner & Biing-Hwang Juang、音声認識の基本(1993)に開示されている。パラメータの組は好ましくはフレーム(周期的フレームに分割される)に基づく。音響プロセッサ 1 0 4 はデジタルシグナルプロセッサ(DSP)として実現される。DSP は音声符号器を含む。一方、音響プロセッサ 1 0 4 は音声符号器として実現される。

【 0 0 1 6 】

パラメータの各フレームは特徴抽出器 1 0 6 に供給される。特徴抽出器 1 0 6 において終了点検出器 1 1 6 は、発話(すなわちワード)の終了点を検出するために抽出されたパラメータを使用する。一実施形態において終了点検出は、米国特許出願第 0 9 / 2 4 6 4 1 4 号(1999年2月8日出願、名称:ノイズの存在下での音声の正確な終了点のための方法及び装置、この米国出願は、本発明の譲受人に譲渡され、参考としてここにその全体が組み込まれている。この出願は米国特許第 6 3 2 4 5 0 9 号として権利化された)に開示されている。この技術に従って、発話は第 1 の開始点及び発話の第 1 の終了点を決定するために、例えば信号対ノイズ比(SNR)しきい値などの第 1 のしきい値と比較される。次に、前記第 1 の開始点に先立つ発話の部分が第 2 の SNR しきい値と比較されて発話の第 2 の開始点が決定される。第 1 の終了点に続く発話の部分は次に第 2 の SNR しきい値と比較されて前記発話の第 2 の終了点が決定される。第 1 及び第 2 の SNR しきい値は好ましくは周期的に再計算され、第 1 の SNR しきい値は好ましくは第 2 の SNR しきい値を越える。

【 0 0 1 7 】

検出された発話に対する周波数領域パラメータのフレームは、時間クラスタリング音声細分化モジュール 1 1 8 に供給されて、一実施形態に従って、米国特許出願第 0 9 / 2 2 5 8 9 1 号(1999年1月4日出願、名称:音声信号の細分化及び認識のためのシステム及び方法、この米国出願は、本発明の譲受人に譲渡され、参考としてここにその全体が組み込まれている。この出願は米国特許第 6 3 2 4 5 0 9 号として権利化された)に記載された圧縮技術を実行する。この技術に従って、周波数領域パラメータにおける各音声フレームは、音声フレームに関連した少なくとも 1 つのスペクトル値によって表現される。次に、隣接フレームの各対に対してスペクトル相異値が決定される。隣接フレームの各対間に初期的クラスタ境界が設定され、当該パラメータにクラスタが生成され、分散値が各クラスタに割り当てられる。分散値は好ましくは、所定のスペクトル相異値の 1 つに等し

10

20

30

40

50

い。次に、複数のクラスタマージパラメータが計算される。クラスタマージパラメータの各々是一对の隣接クラスタに関連する。最小クラスタマージパラメータは複数のクラスタマージパラメータから選択される。マージされたパラメータは次に、最小クラスタマージパラメータに関連したクラスタ間のクラスタ境界を相殺して、マージされた分散値をマージされたクラスタに割り当てることによって形成される。マージされた分散値は最小クラスタマージパラメータに関連したクラスタに割り当てられた分散値を表す。このプロセスは好ましくは複数のマージされたクラスタを形成するために反復され、好ましくは複数のマージされたクラスタに従って分割された音声信号が形成される。

【 0 0 1 8 】

時間クラスタリング音声細分化モジュール 1 1 8 は、例えば時間正規化モジュールなどの他の装置により置き換えられる。しかしながら、時間クラスタリング音声細分化モジュール 1 1 8 は、以前のフレームと比較したときに最小の相異をもつフレームをクラスタにマージし、個々のフレームの代わりに算術平均 (mean average) を使用するので、時間クラスタリング音声細分化モジュール 1 1 8 は、処理された発話におけるより多くの情報を使用する。時間クラスタリング音声細分化モジュール 1 1 8 は好ましくは、当業界で知られており以下に説明する動的時間ゆがみ (D T W) モデルを使用するパターン比較ロジック 1 1 0 に関連して使用される。

【 0 0 1 9 】

クラスタ手段は音声レベル正規化器 1 2 0 に供給される。一実施形態において、音声レベル正規化器 1 2 0 は、各クラスタにチャンネルあたり平均 2 ビット (すなわち、周波数ごとに 2 つのビット) を割り当てることによって音声振幅を量子化する。ケプストラム係数が抽出される他の実施形態において、音声レベル正規化器 1 2 0 は、当業者により理解されるように、クラスタ手段を量子化するのに使用されない。音声レベル正規化器 1 2 0 により生成された出力は特徴抽出器 1 0 6 によってパターン比較ロジック 1 1 0 に供給される。

【 0 0 2 0 】

音声認識サブシステム 1 0 0 のボキャブラリワードのすべてに対するテンプレートの組は、テンプレートデータベース 1 0 8 内に永久的に記憶される。テンプレートの組は好ましくは、以下に説明するような話者独立テンプレート構築サブシステムにより構築される一組の話者独立テンプレートである。テンプレートデータベース 1 0 8 は好ましくは、例えばフラッシュメモリなどの任意の従来の形態の不揮発性記憶媒体として実現される。このことは、音声認識サブシステム 1 0 0 に対する電源が O F F されたときにテンプレートがテンプレートデータベース 1 0 8 内に残ることを可能にする。

【 0 0 2 1 】

パターン比較ロジック 1 1 0 は、特徴抽出器 1 0 6 からのベクトルをテンプレートデータベース 1 0 8 内に記憶されたすべてのテンプレートと比較する。比較結果すなわちベクトルとテンプレートデータベース 1 0 8 内に記憶されたすべてのテンプレート間の距離は、決定ロジック 1 1 2 に供給される。決定ロジック 1 1 2 は、所定のマッチングしきい値内で N 個の最も近いマッチ (合致) を選択する従来の “N - ベスト” 選択アルゴリズムを使用する。ユーザは次にどちらの選択が意図されたかについて質問される。決定ロジック 1 1 2 の出力はボキャブラリ内のどのワードが話されたかに関する決定である。

【 0 0 2 2 】

一実施形態において、パターン比較ロジック 1 1 0 及び決定ロジック 1 1 2 はコンバージェンスを試験するために D T W 技術を使用する。D T W 技術は当業界で知られており、その全体がここに参照として組み込まれた、Lawrence Rabiner & Biing-Hwang Juang 音声認識の基本 200-238(1993)に記載されている。D T W 技術に従って、試験すべき発話の時間シーケンスをテンプレートデータベース 1 0 8 内に記憶された各発話に対する時間シーケンスに対してプロットすることによって、トレリスが形成される。次に、試験されている発話は、ポイント (例えば各 1 0 m s) ごとに一度に 1 発話だけテンプレートデータベース 1 0 8 内の各発話と比較される。テンプレートデータベース 1 0 8 における各発

10

20

30

40

50

話に対して、試験されている発話は時間上で調整すなわち“ゆがみが加えられ”、テンプレートデータベース108における発話に最も近い可能なマッチが達成されるまで、特定の点で圧縮あるいは伸張される。時間の各点で2つの発話が比較されて、マッチがその点（零コスト）で宣言されるかあるいは、ミスマッチが宣言される。特定の点でのミスマッチがあった場合には、試験されるべき発話が圧縮、伸張あるいは必要に応じてミスマッチされる。当該プロセスは2つの発話が互いに完全に比較されるまで反復される。多数（概して数千）の別々に調整された発話が可能である。最も低い（すなわち、最も少ない数の圧縮及び/又は拡大及び/又はミスマッチを要求する）コスト関数が選択される。ビタビ復号アルゴリズムの場合と同様に、最も低い全コストをもつ経路を決定するために、好ましくはテンプレートデータベース108における発話の各ポイントから後方を見ることによって選択が実行される。このことは、個別に調整された発話のすべての可能性を生成する“強引な”方法に頼ることなしに最も低いコスト（すなわち最も近いマッチ）の調整された発話を決定することを可能にする。次に、テンプレートデータベース108における発話のすべてに対する最も低いコストの調整された発話は次に比較されて、最も低いコストをもつ1つが、試験された発話に最も近くマッチした記憶された発話として選択される。

10

【0023】

パターン比較ロジック110及び決定ロジック112は好ましくはマイクロプロセッサとして実現される。音声認識サブシステム100は例えばASICである。音声認識サブシステム100の認識精度は、音声認識サブシステム100がどれくらい正確にボキャブラリ内の話されたワードまたはフレーズを認識するかについての測定基準となる。例えば、95%の認識精度は、音声認識サブシステム100がボキャブラリ内のワードを100のうち95回正確に認識することを示す。

20

【0024】

一実施形態によれば、音声認識サブシステム（図示せず）は、音声認識サブシステムに対する音声入力を認識するために図3のフローチャートに示されたアルゴリズムステップを実行する。ステップ200において、入力音声は音声認識サブシステムに供給される。次に制御フローはステップ202に進む。ステップ202において発話の終了点が検出される。特定の実施形態において、発話の終了点は図2を参照して上記したように上記の米国特許出願第09/246414号に記載された技術に従って検出される。制御フローは次にステップ204に進む。

30

【0025】

ステップ204において時間クラスタリング音声細分化が抽出された発話に関して実行される。特定の実施形態において、使用される時間クラスタリング音声細分化技術は、図2を参照して上記したように、上記した米国特許出願第09/225891号（米国特許第6278972号として権利化された）に記載された技術である。制御フローは次にステップ208に進む。ステップ206においてステップ204において生成された音声クラスタ手段とのマッチのために話者独立テンプレートが提供される。話者独立テンプレートは好ましくは図4-6を参照して以下に述べる技術に従って構成される。制御フローは次にステップ208に進む。ステップ208において特定の発話に対するクラスタとすべての話者独立テンプレート間でDTWマッチが実行され、最も近いマッチングテンプレートが認識された発話として選択される。特定の実施形態においてDTWマッチがLawrence Rabiner & Biing-Hwang Juang 音声認識の基本200-238(1993)に記載され、図2を参照して上記した技術に従って実行される。時間クラスタリング音声細分化以外の方法がステップ204において実行される。そのような方法は、例えば時間正規化を含む。

40

【0026】

一実施形態によれば、図4に示すように、話者独立テンプレート構築サブシステム300は、プロセッサ302と記憶媒体304を含む。プロセッサ100は好ましくはマイクロプロセッサであるが、他の従来の形態のプロセッサ、専用プロセッサ、デジタルシグナルプロセッサ(DSP)、コントローラ、あるいはステートマシンを含む。プロセッサ3

50

02は、フラッシュメモリ、EEPROMメモリ、RAMメモリ、ファームウェア命令を保持するように構成されたROMメモリ、プロセッサ302上で動作するソフトウェアモジュール、あるいは任意の他の従来の形態のメモリとして実現される記憶媒体304に結合される。話者独立テンプレート構築サブシステム300は好ましくは、ユニックスオペレーティングシステム上で動作するコンピュータとして実現される。他の実施形態において、記憶媒体304はオンボードRAMメモリあるいはプロセッサ302であり、記憶媒体304はASIC内に存在する。一実施形態において、プロセッサ302は、図6を参照して以下に述べるステップなどのアルゴリズムステップを実行するために記憶媒体304により記憶される一連の命令を実行するように構成される。

【0027】

他の実施形態によれば、図5に示すように、話者独立テンプレート構築サブシステム400は、終了点検出器402、時間クラスタリング音声細分化ロジック404、ベクトル量子化器406、コンバージェンス試験器408、そしてK-平均音声細分化ロジック410を含む。制御プロセッサ(図示せず)は好ましくは、話者独立テンプレート構築サブシステム400が実行する反復の回数を制御するために使用される。

【0028】

終了点検出器402は時間クラスタリング音声細分化ロジック404に接続される。時間クラスタリング音声細分化ロジック404は、ベクトル量子化器406に結合される。ベクトル量子化器406はコンバージェンス試験器408及びK-平均音声細分化ロジック410に結合される。制御プロセッサは好ましくは、制御バス(図示せぬ)を介して終了点検出器402、時間クラスタリング音声細分化ロジック404、ベクトル量子化器406、コンバージェンス試験器408、そしてK-平均音声細分化ロジック410に結合される。

【0029】

学習すべき発話の学習サンプル $S_x(n)$ はフレームで終了点検出器402に供給される。学習サンプルは好ましくは、学習すべき発話記憶されるトレーニングデータベース(図示せぬ)から供給される。一実施形態において、トレーニングデータベースは100のワードを含み、各ワードは、全部で10000の記憶された発話に対して100の異なる話者により話される。終了点検出器402は発話の開始点と終了点とを検出する。一実施形態において、終了点検出器402は、上記した米国特許出願第09/246414号及び図2を参照して上記した技術に従って動作する。

【0030】

終了点検出器402は、検出された発話を時間クラスタリング音声細分化ロジック404に供給する。時間クラスタリング音声細分化ロジック404は検出された発話に関して圧縮アルゴリズムを実行する。一実施形態において、時間クラスタリング音声細分化ロジック404は、上記した米国特許出願第09/225891号及び図2を参照して上記した技術に従って動作する。一実施形態において、時間クラスタリング音声細分化ロジック404は検出された発話を20のセグメントに分割するが、各セグメントはクラスタ手段を含む。

【0031】

時間クラスタリング音声細分化ロジック404は、所定のワードに対するトレーニング用発話のすべてに対するクラスタ手段をベクトル量子化器406に供給する。ベクトル量子化器406は、発話に対する(すなわち同じワードのすべての話者に対する)クラスタ手段をベクトル量子化し、最終的なベクトルを当該発話に対する潜在的な話者独立(SI)テンプレートとしてコンバージェンス試験器408に供給する。ベクトル量子化器406は好ましくは、種々の既知のベクトル量子化(VQ)技術に従って動作する。種々のVQ技術は、例えば、Gersho & R.M. Gray, ベクトル量子化及び信号圧縮(1992)に記載されている。特定の実施形態において、ベクトル量子化器406は、4つのクラスタベクトルを生成する。すなわち、例えば、各セグメントは、各セグメントを4つのクラスタとして表すベクトル量子化器406に直列に供給される。各クラスタは特定のワードに対す

10

20

30

40

50

る各話者を表わし、ワードあたり複数のクラスタが存在する。一実施形態によれば、テンプレートあたり80(4つのクラスタ×20セグメント)のベクトルが存在する。

【0032】

コンバージェンス試験器408は、潜在的なS Iテンプレートと試験すべき発話の試験サンプル $S_y(n)$ とを比較する。試験サンプルはフレームの形態でコンバージェンス試験器408に供給される。試験用サンプルは好ましくは、試験すべき発話が記憶される試験データベース(図示せぬ)から供給される。一実施形態において、試験用データベースは100のワードを含み、各ワードは全部で10000の記憶された発話に対して100の異なる話者により話される。ワードは好ましくは、トレーニングデータベースに含まれている同じワードであるが、100の異なる話者によって話される。コンバージェンス試験器408は、訓練すべき発話に対する潜在的なS Iテンプレートと試験すべき発話に対するサンプルとを比較する。一実施形態において、コンバージェンス試験器408は、コンバージェンスを試験するためのDTWアルゴリズムを使用するように構成される。使用されるDTWアルゴリズムは好ましくは、Lawrence Rabiner & Biing-Hwang Juang 音声認識の基本200-238(1993)に記載され、図2を参照して上記した技術である。

【0033】

一実施形態において、コンバージェンス試験器408は、前記データベース内のすべてのワードに対する結果の精度と、前記データベースの分散とを潜在的なS Iテンプレートで解析するように構成される。分散がまずチェックされ、この分散が所定のしきい値以下であるならば、精度がチェックされる。分散は好ましくはセグメントごとに計算され、次に加算されて全体の分散値が取得される。特定の実施形態において、分散は4つのクラスタのベストマッチ(最善の合致)に対する平均二乗誤差を計算することによって取得される。平均二乗誤差の技術は、当業界で良く知られている。試験データベースからの発話がトレーニングデータベースにより生成された潜在的なS Iテンプレートと一致する(すなわち認識がデータベース内のすべてのワードに対して正しい)ならば、正確であるとみなされる。

【0034】

潜在的なS Iテンプレートは、ベクトル量子化器406からK-平均音声細分化ロジック410に供給される。K-平均音声細分化ロジック410は好ましくはフレームに分割されたトレーニングサンプルを受信する。コンバージェンス試験器408がコンバージェンスに対する第1の試験を実行した後に、分散または精度に対する結果は、分散及び精度に対する所定のしきい値以下になる。一実施形態において、分散または精度に対する結果が所定のしきい値以下であるならば、他の反復が実行される。したがって、制御プロセッサは、K-平均音声細分化ロジック410に対してトレーニングサンプルに関してK-平均細分化を実行するように命令し、以下に述べるような分割された音声フレームを生成する。K-平均音声細分化に従って、トレーニングサンプルは好ましくはDTW技術を用いて潜在的なS Iテンプレートとマッチングがとられて、図2を参照して上記したような最適なパスを生成する。次にトレーニングサンプルは最適なパスに従って分割される。例えば、トレーニングサンプルの第1の5つのフレームは、潜在的なS Iテンプレートの第1のフレームとマッチし、トレーニングサンプルの次の3つのフレームは潜在的なS Iテンプレートの第2のフレームにマッチし、トレーニングサンプルの次の10のフレームは、潜在的なS Iテンプレートの第3のフレームにマッチする。この場合、トレーニングサンプルの第1の1つのフレームは1つのフレームに分割され、次の3つのフレームは第2のフレームに分割され、次の10のフレームは第3のフレームに分割される。一実施形態において、K-平均音声細分化ロジック410は、Lawrence Rabiner & Biing-Hwang Juang 音声認識の基本382-384(1993)(この文献の全体が参照としてここに組み込まれている)に記載された例示的なK-平均細分化技術に従ってK-平均細分化を実行する。次にK-平均音声細分化ロジック410はクラスタ手段の更新されたフレームをベクトル量子化器406に供給する。ベクトル量子化器406は、クラスタ手段をベクトル量子化して(新たな潜在的なS Iテンプレートを具備する)最終的なベクトルをコンバージェンス試験器

10

20

30

40

50

408に供給して他のコンバージェンス試験を実行する。当業者ならば、この反復工程が所定のしきい値以上の分散及び精度の結果を得るのに必要なだけ継続されることを認識する。

【0035】

コンバージェンス試験がパスすると、好ましくは潜在的な（ここでは最終的な）S I テンプレートが図2の音声認識サブシステムなどの音声認識サブシステムにおいて使用される。最終的なS I テンプレートは、図2のテンプレートデータベース108に記憶されるかあるいは図3のフローチャートのステップ206において使用される。

【0036】

一実施形態において、話者独立テンプレート構築サブシステム（図示せず）は、発話に対する話者独立テンプレートを構築するために図6のフローチャートにおいて例示された方法ステップを実行する。ステップ500において発話のトレーニングサンプルがトレーニングデータベース（図示せぬ）から入手される。トレーニングデータベースは好ましくは、多数のワード（例えば100ワード）を含む。各ワードは多数の話者（ワードあたり例えば100の話者）によって話される。制御フローは次にステップ502に進む。

10

【0037】

ステップ502において、発話を検出するために終了点検出がトレーニングサンプルに関して実行される。一実施形態において終了点検出は、上記した米国特許出願第09/246414号及び図2を参照して上記した技術に従って動作する。制御フローは次にステップ504に進む。

20

【0038】

ステップ504において、検出された発話に関して時間クラスタリング音声細分化が実行され、当該発話を複数のセグメントに圧縮する。各セグメントは手段によって表示される。特定の実施形態において、発話は20のセグメントに分割される。各セグメントはクラスタ手段を含む。一実施形態において、発話は20のセグメントに圧縮される。各セグメントはクラスタ手段を含む。一実施形態において時間クラスタリング音声細分化は、上記した米国特許出願第09/225891号（この出願は米国特許第6278972号として権利化された）及び図2を参照して上記した技術に従って動作される。制御フローは次にステップ506に進む。

30

【0039】

ステップ506において、同じワードのすべての話者に対するトレーニングサンプルに対するクラスタ手段は、ベクトル量子化される。特定の実施形態において、クラスタ手段は、A.Gersho & R.M. Gray, ベクトル量子化及び信号圧縮(1992)に記載された種々の知られたVQ技術の1つに従ってベクトル量子化される。特定の実施形態において、4つのクラスタベクトルが生成される。すなわち、例えば、各セグメントは4つのクラスタとして表される。各クラスタは特定のワードに対する各話者を表し、ワードごとに複数のクラスタが存在する。一実施形態によれば、テンプレートあたり80のベクトル（4つのクラスタ×20セグメント）が生成される。制御フローは次にステップ510に進む。

【0040】

ステップ508において試験データベース（図示せぬ）から入手した試験サンプルは、コンバージェンスが試験される。試験データベースは好ましくはトレーニングデータベースに含まれる同じワードを含む。各ワードは多数の話者（発話あたり例えば100の話者）によって話される。

40

【0041】

ステップ510において、量子化されたベクトルは、コンバージェンスを試験するために潜在的なS I テンプレートとして試験サンプルと比較される。一実施形態において、コンバージェンス試験はD T Wアルゴリズムである。使用されるD T Wアルゴリズムは、Lawrence Rabiner & Biing-Hwang Juang、音声認識の基本,200-238（1993）及び図2を参照して上記した技術である。

【0042】

50

一実施形態において、ステップ510のコンバージェンス試験では、データベース内のすべてのワードに対する結果の精度と、潜在的なS Iサンプルを有するデータベースの分散である。分散がまずチェックされ、この分散が所定のしきい値以下ならば、次に精度がチェックされる。分散は好ましくは、セグメントごとに計算され、次に加算されて全体の分散値が取得される。特定の実施形態において、分散は、4つのクラスタの最善の合致に対する平均二乗誤差を計算することによって取得される。平均二乗誤差の技術は当業界でよく知られている。試験データベースによって生成された潜在的なS Iテンプレートがトレーニングデータベースからの発話に合致するならば(すなわち、データベース内のすべてのワードに対して認識が正しいならば)、コンバージェンス試験は正確であるとみなされる。制御フローは次にステップ512に進む。

10

【0043】

ステップ512において、分散あるいは精度に対するステップ510のコンバージェンス試験の結果が、分散及び精度に対する所定のしきい値以下ならば、他の反復が実行される。したがって、トレーニングサンプルに関してK - 平均音声細分化が実行される。K - 平均音声細分化はトレーニングサンプルを好ましくはDTW技術を用いて潜在的なS Iテンプレートとマッチさせ、図2を参照して上記したような最適なパスを生成する。トレーニングサンプルは次に最適なパスに従って分割される。一実施形態において、K - 平均音声細分化は、Lawrence Rabiner & Biing-Hwang Juang、音声認識の基本,382-384(1993)に記載された技術に従って実行される。制御フローは次にステップ506に戻って、クラスタ手段の更新されたフレームがベクトル量子化され、ステップ510において、試験データベースからのサンプルによってコンバージェンスが(新たな潜在的なS Iテンプレートとして)試験される。当業者ならば、この反復工程は予め決められたしきい値以上の分散及び精度結果を達成するのに必要なだけ継続されることを認識する。

20

【0044】

いったんコンバージェンス試験がパスすると(すなわち、いったんしきい値が達成されると)、潜在的な(ここでは最終的な)S Iテンプレートは好ましくは、図2の音声認識サブシステムなどの音声認識サブシステムにおいて使用される。最終的なS Iテンプレートは、図2のテンプレートデータベース108に記憶されるかあるいは図3のフローチャートのステップ206において使用される。

【0045】

すなわち、話者独立音声認識システムに対する音声テンプレートを構築するための新規かつ改善された方法及び装置が説明された。当業者ならば、データ、命令、コマンド、情報、信号、ビット、シンボルそして上記記述において参照されるチップは好ましくは、電圧、電流、電磁波、磁界あるいは粒子、オプティカルフィールド、またはそれらの任意の組み合わせであることを認識するであろう。当業者ならば、種々の例示的な論理ブロック、モジュール、回路そしてここで開示された実施形態に関連したアルゴリズムステップは、電子ハードウェア、コンピュータソフトウェアあるいは両方の組み合わせとして実現される。種々の例示的な要素、ブロック、モジュール、回路、そしてステップは概してそれらの機能の観点から記載された。機能がハードウェアとして実現されるかソフトウェアとして実現されるかは特定のアプリケーションと全体のシステムに課される設計上の拘束に依存する。熟練した技術者ならば、このような環境の下でハードウェアとソフトウェアとを交換したり、記述された機能を各特定のアプリケーションに対して最適に実行する方法を認識するであろう。例として、種々の例示的なロジカルブロック、モジュール、回路そしてここに開示された実施形態に関連して記述されたアルゴリズムステップは、デジタルシグナルプロセッサ(DSP)、特定用途向け集積回路(ASIC)、フィールドプログラマブルゲートアレイ(FPGA)あるいは他のプログラマブルロジックデバイス、ディスクリットゲートまたはトランスファゲート、例えばレジスタやFIFOなどのディスクリットハードウェア要素、一連のファームウェア命令を実行するプロセッサ、任意の従来のプログラマブルソフトウェアモジュール及びプロセッサ、あるいはここに記載された機能を実行するために設計されたそれらの任意の組み合わせ、である。プロセッサは好ま

30

40

50

しくはマイクロプロセッサであるが、変形例として、プロセッサは任意の従来のプロセッサ、コントローラ、マイクロコントローラ、あるいは他のマシンである。ソフトウェアモジュールは、RAMメモリ、フラッシュメモリ、ROMメモリ、EEPROMメモリ、レジスタ、ハードウェアディスク、リムーバブルディスク、CD-ROM、あるいは任意の他の形態の公知の記憶媒体である。例示的なプロセッサは好ましくは記憶媒体に結合され、当該記憶媒体から情報を読み出したり、当該記録媒体に情報を書き込む。あるいは、記憶媒体はプロセッサと一体である。プロセッサ及び記憶媒体はASIC内に配置される。ASICは電話内に配置される。あるいは、プロセッサ及び記憶媒体が電話内に配置される。プロセッサはDSPとマイクロプロセッサとの組み合わせとしてあるいはDSPコアなどに関連する2つのマイクロプロセッサとして実現される。

10

【0046】

本発明の好ましい実施形態が示されかつ記述された。しかしながら、本発明の精神あるいは範囲から逸脱することなしにここに開示された実施形態に対して種々の変形を行うことができる。したがって、本発明は次の特許請求の範囲に従う以外には限定されることはない。

【図面の簡単な説明】

【図1】 図1は、話者独立音声認識に対する音声テンプレートを構築して実行するためのシステムのブロック図である。

【図2】 図2は、図1のシステムにおいて使用可能な音声認識サブシステムのブロック図である。

20

【図3】 図3は、入力音声サンプルを認識するために、図2のサブシステムなどの、音声認識サブシステムによって実行される方法ステップを示すフローチャートである。

【図4】 図4は、図1のシステムにおいて使用可能なテンプレート構築サブシステムのブロック図である。

【図5】 図5は、図1のシステムにおいて使用可能なテンプレート構築サブシステムのブロック図である。

【図6】 図6は、音声テンプレートを構築するために、図4のサブシステムや図5のサブシステムなどの、テンプレート構築サブシステムによって実行される方法ステップを示すフローチャートである。

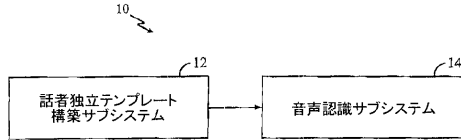
【符号の説明】

30

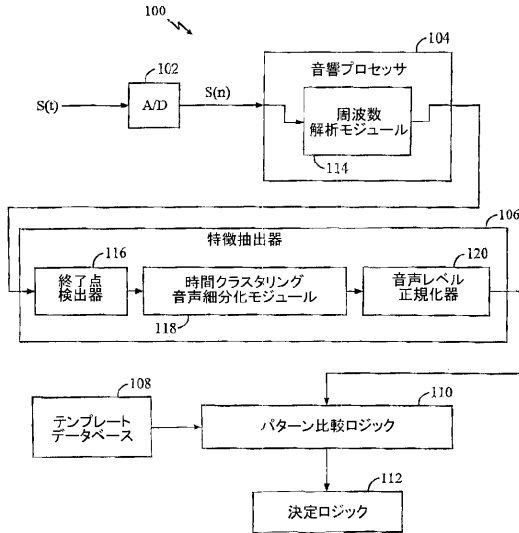
- 10 システム
- 12 話者独立テンプレート構築サブシステム
- 14 音声認識サブシステム
- 100 音声認識サブシステム
- 102 アナログ/デジタルコンバータ(A/D)
- 104 フロントエンドプロセッサ
- 106 特徴抽出器
- 108 音声テンプレートデータベース
- 110 パターン比較ロジック
- 112 決定ロジック
- 114 周波数解析モジュール
- 116 エンドポイントモジュール
- 118 時間クラスタリング音声細分化モジュール
- 120 音声レベル正規化器

40

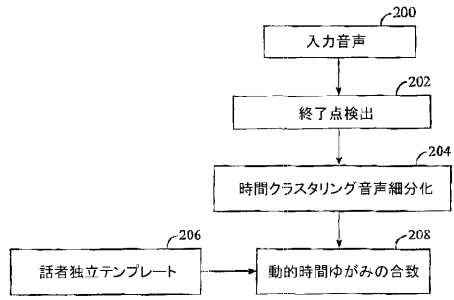
【図1】



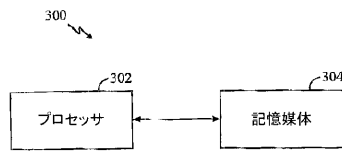
【図2】



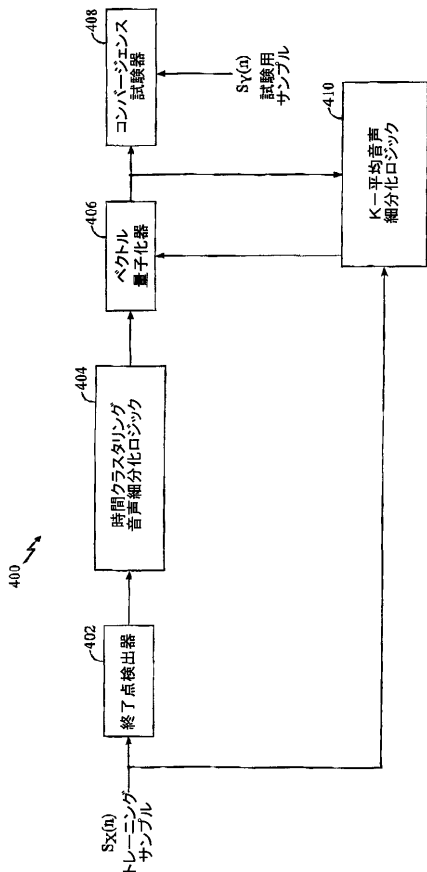
【図3】



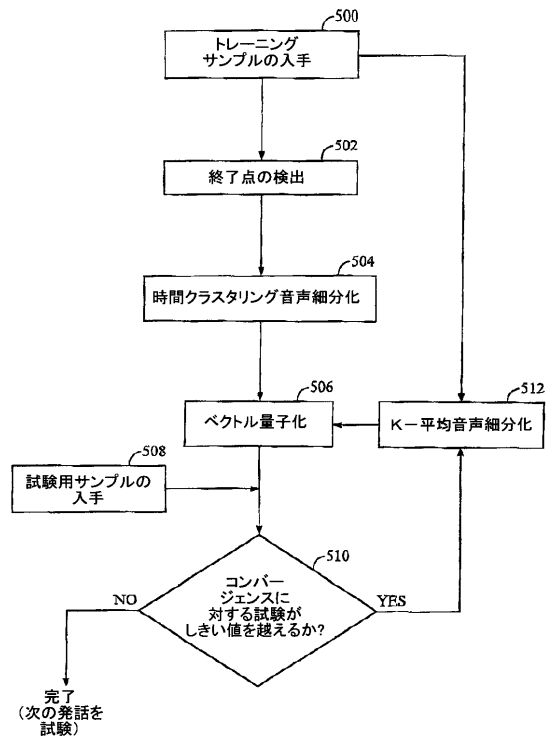
【図4】



【図5】



【図6】



フロントページの続き

(74)代理人 100084618

弁理士 村松 貞男

(74)代理人 100092196

弁理士 橋本 良郎

(72)発明者 ビー、ニン

アメリカ合衆国、カリフォルニア州 9 2 1 2 8 サン・ディエゴ、ブリーズウェイ・プレイス
1 4 2 0 9

審査官 渡邊 聡

(56)参考文献 特開昭62-217292(JP,A)

特開昭62-072000(JP,A)

特許第2749811(JP,B2)

特開平07-512966(JP,A)

特開平05-241591(JP,A)

特開平11-212592(JP,A)

(58)調査した分野(Int.Cl., DB名)

G10L 15/06