

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4061094号  
(P4061094)

(45) 発行日 平成20年3月12日(2008.3.12)

(24) 登録日 平成19年12月28日(2007.12.28)

(51) Int.Cl.			F I		
<b>G 1 0 L</b>	<b>15/14</b>	<b>(2006.01)</b>	G 1 0 L	15/14	2 0 0 C
<b>G 1 0 L</b>	<b>15/06</b>	<b>(2006.01)</b>	G 1 0 L	15/06	3 1 0 T
<b>G 1 0 L</b>	<b>15/20</b>	<b>(2006.01)</b>	G 1 0 L	15/20	2 0 0 Q

請求項の数 13 (全 13 頁)

(21) 出願番号	特願2002-72456 (P2002-72456)	(73) 特許権者	390009531
(22) 出願日	平成14年3月15日(2002.3.15)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公開番号	特開2003-280686 (P2003-280686A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公開日	平成15年10月2日(2003.10.2)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
審査請求日	平成14年12月20日(2002.12.20)	(74) 代理人	100086243
審判番号	不服2005-21439 (P2005-21439/J1)		弁理士 坂口 博
審判請求日	平成17年11月4日(2005.11.4)	(74) 代理人	100091568
			弁理士 市位 嘉宏

最終頁に続く

(54) 【発明の名称】 音声認識装置、その音声認識方法及びプログラム

(57) 【特許請求の範囲】

【請求項1】

所定の音声と、予め収録された音声データの音韻隠れマルコフモデルとのマッチングを取ることにより音声認識を行う音声認識装置において、

認識対象である入力音声の特徴量を抽出する特徴量抽出部と、

前記音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成される雑音源ごとに独立の隠れマルコフモデルとを合成し雑音源ごとに合成モデルを作成する合成モデル作成部と、

前記特徴量抽出部にて抽出された前記入力音声の特徴量と前記合成モデル作成部にて作成された雑音源ごとの前記合成モデルとのマッチングを取ることにより前記入力音声

10

認識する音声認識部とを備え、  
前記音声認識部は、前記入力音声における発話区間を区切る適当な区間ごとに当該入力音声の特徴量と前記合成モデルとの尤度計算を行ってマッチングを取ることを特徴とする音声認識装置。

【請求項2】

前記音声認識部は、前記入力音声の音声フレームごとに、独立してマッチング対象となる前記合成モデルを選択し、当該入力音声の特徴量と当該合成モデルとの尤度計算を行ってマッチングを取ることを特徴とする請求項1に記載の音声認識装置。

【請求項3】

所定の音声と、予め収録された音声データの音韻隠れマルコフモデルとのマッチングを

20

取ることにより音声認識を行う音声認識装置において、

認識対象である入力音声の特徴量を抽出する特徴量抽出部と、

前記音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成される雑音源ごとに独立の隠れマルコフモデルとを合成し雑音源ごとに合成モデルを作成する合成モデル作成部と、

前記特徴量抽出部にて抽出された前記入力音声の特徴量と前記合成モデル作成部にて作成された雑音源ごとの前記合成モデルとのマッチングを取ることにより前記入力音声を認識する音声認識部とを備え、

前記音声認識部は、発話中の前記入力音声に含まれる雑音の変化に応じてマッチング対象となる前記合成モデルを動的に選択しながら当該入力音声の特徴量との尤度計算を行ってマッチングを取ることを特徴とする音声認識装置。

10

【請求項 4】

音声認識のためのモデルとなる音声データを格納した音声データベースと、

所定の雑音環境で発生が想定される雑音データを格納した雑音データベースと、

前記音声データベースから読み出した音声データに基づき作成される音声モデルと前記雑音データベースから読み出した雑音データに基づき雑音源ごとに独立に作成される雑音モデルとを合成して雑音源ごとに合成モデルを作成する合成モデル作成部と、

認識対象である入力音声の特徴量と前記合成モデル作成部にて作成された雑音源ごとの前記合成モデルとに関して、前記入力音声の音声フレームごとに独立して尤度計算を行ってマッチングを取ることにより音声認識を行う音声認識部と

20

を備えることを特徴とする音声認識装置。

【請求項 5】

音声認識のためのモデルとなる音声データを格納した音声データベースと、

所定の雑音環境で発生が想定される雑音データを格納した雑音データベースと、

前記音声データベースから読み出した音声データに基づき作成される音声モデルと前記雑音データベースから読み出した雑音データに基づき雑音源ごとに独立に作成される雑音モデルとを合成して雑音源ごとに合成モデルを作成する合成モデル作成部と、

認識対象である入力音声の特徴量と前記合成モデル作成部にて作成された雑音源ごとの前記合成モデルとに関して、発話中の前記入力音声に含まれる雑音の変化に応じてマッチング対象となる前記合成モデルを動的に選択しながら当該入力音声の特徴量との尤度計算を行ってマッチングを取ることにより音声認識を行う音声認識部と

30

を備えることを特徴とする音声認識装置。

【請求項 6】

コンピュータを制御して、音声を認識する音声認識方法において、

認識対象である入力音声の特徴量を抽出し、メモリに格納するステップと、

所定の音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成される雑音源ごとに独立の隠れマルコフモデルとを合成して生成された雑音源ごとに独立の合成モデルをメモリから読み出すステップと、

前記入力音声の音声フレームごとに、前記メモリに格納された前記入力音声の特徴量と雑音源ごとの前記合成モデルとの尤度計算を行ってマッチングを取り、当該マッチングの結果に基づいて前記入力音声を認識するステップと

40

を含むことを特徴とする音声認識方法。

【請求項 7】

前記入力音声を認識するステップは、前記入力音声の音声フレームごとに、独立してマッチング対象となる前記合成モデルを選択し、当該入力音声の特徴量と当該合成モデルとの尤度計算を行ってマッチングを取ることを特徴とする請求項 6 に記載の音声認識方法。

【請求項 8】

コンピュータを制御して、音声を認識する音声認識方法において、

認識対象である入力音声の特徴量を抽出し、メモリに格納するステップと、

所定の音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成され

50

る雑音源ごとに独立の隠れマルコフモデルとを合成して生成された雑音源ごとに独立の合成モデルをメモリから読み出すステップと、

前記メモリに格納された前記入力音声の特徴量と雑音源ごとの前記音韻隠れマルコフモデルとに関して、発話中の前記入力音声に含まれる雑音の変化に応じてマッチング対象となる前記合成モデルを動的に選択しながら当該入力音声の特徴量との尤度計算を行ってマッチングを取ることにより前記入力音声を認識するステップとを含むことを特徴とする音声認識方法。

【請求項 9】

コンピュータを制御して、音声認識処理を実行させるプログラムであって、  
認識対象である入力音声の特徴量を抽出する特徴量抽出手段と、

予め収録された音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成される雑音源ごとに独立の隠れマルコフモデルとを合成し雑音源ごとに合成モデルを作成する合成モデル作成手段と、

前記入力音声における発話区間を区切る適当な区間ごとに、前記入力音声の特徴量と雑音源ごとの前記合成モデルとの尤度計算を行ってマッチングを取ることにより前記入力音声を認識する音声認識手段として、

前記コンピュータを機能させることを特徴とするプログラム。

【請求項 10】

前記プログラムによる前記音声認識手段は、前記入力音声の音声フレームごとに、独立してマッチング対象となる前記合成モデルを選択し、当該入力音声の特徴量と当該合成モデルとの尤度計算を行ってマッチングを取ることとする請求項 9 に記載のプログラム。

【請求項 11】

コンピュータを制御して、音声認識処理を実行させるプログラムにおいて、  
認識対象である入力音声の特徴量を抽出する特徴量抽出手段と、

予め収録された音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成される雑音源ごとに独立の隠れマルコフモデルとを合成し雑音源ごとに合成モデルを作成する合成モデル作成手段と、

発話中の前記入力音声に含まれる雑音の変化に応じてマッチング対象となる前記合成モデルを動的に選択しながら、前記入力音声の特徴量と雑音源ごとの前記合成モデルとの尤度計算を行ってマッチングを取ることにより前記入力音声を認識する音声認識手段として

、  
前記コンピュータを機能させることを特徴とするプログラム。

【請求項 12】

コンピュータを制御して音声認識処理を実行させるプログラムを、当該コンピュータが読み取り可能に記録した記録媒体であって、

前記プログラムは、  
認識対象である入力音声の特徴量を抽出する特徴量抽出手段と、

予め収録された音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成される雑音源ごとに独立の隠れマルコフモデルとを合成し雑音源ごとに合成モデルを作成する合成モデル作成手段と、

前記入力音声における音声フレームごとに、前記入力音声の特徴量と雑音源ごとの前記合成モデルとの尤度計算を行ってマッチングを取ることにより前記入力音声を認識する音声認識手段として、

前記コンピュータを機能させることを特徴とする記録媒体。

【請求項 13】

コンピュータを制御して音声認識処理を実行させるプログラムを、当該コンピュータが読み取り可能に記録した記録媒体であって、

前記プログラムは、  
認識対象である入力音声の特徴量を抽出する特徴量抽出手段と、

10

20

30

40

50

予め収録された音声データの音韻隠れマルコフモデルと予め収録された雑音データから生成される雑音源ごとに独立の隠れマルコフモデルとを合成し雑音源ごとに合成モデルを作成する合成モデル作成手段と、

発話中の前記入力音声に含まれる雑音の変化に応じてマッチング対象となる前記合成モデルを動的に選択しながら、前記入力音声の特徴量と雑音源ごとの前記合成モデルとの尤度計算を行ってマッチングを取ることにより前記入力音声を認識する音声認識手段として

、前記コンピュータを機能させることを特徴とする記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、音声認識システムに関し、特に突発的に発生したり不規則に発生したりするような急激な変化を伴う雑音に対処して音声認識を行う方法に関する。

【0002】

【従来の技術】

コンピュータにて音声を認識する音声認識処理では、種々の雑音源が存在する環境下においても精度の高い認識を行うことが課題の1つとなっている。

従来、雑音環境下で音声認識を行う手法として、スペクトラル・サブトラクション (Spectral Subtraction) 法、HMM (Hidden Markov Model: 隠れマルコフモデル) 合成法、CDCN (Codeword-Dependent Cepstral Normalization) 法など、種々の手法が提案されている。

【0003】

これらの手法は、音声認識を目的としていることに鑑み、基本的に、1回の発話(発生)が終了した後に、当該発話中の音声信号の中から雑音に相当する部分を特定し、この特定された雑音部分を考慮して(もしくは除去して)音声認識を行っている。

例えば、HMM合成法では、種々の雑音HMMと音声HMMをそれぞれ合成して雑音の要素が加味された音韻隠れマルコフモデル(合成HMM)を生成し、認識対象の音声に対して最も尤度の高い合成HMMに基づいて音声認識を行うことにより雑音に対処しているが、従来のHMM合成法は、1回の発話ごとに尤度が最大である合成HMMを選択し、認識結果として採用している。すなわち、1つの発話に対して1つの雑音HMMが選択されることとなる。

【0004】

【発明が解決しようとする課題】

ところで、種々の雑音源が存在する環境下では、定常的に発生し続ける雑音や突発的に発生する雑音、不規則に発生する雑音など、雑音の発生の仕方も様々である。上述した従来の音声認識処理における雑音への対処技術は、1回の発話ごとに雑音の種類を認定して対処しているため、定常的に発生し続ける雑音や規則的に発生する雑音に対しては十分な効果を奏し、良好な音声認識を実現することができる。

しかし、突発的に発生する雑音や不規則に発生する雑音は、発話の最中に発生する場合があり、1回の発話ごとに雑音の種類を認定する従来の技術は、このような急激に変化する雑音に対処できず、音声認識の精度を低下させる原因となっていた。

【0005】

そこで、本発明は、突発的に発生する雑音や不規則に発生する雑音などのような急激な変化を伴う雑音に対しても十分に対処し、精度の高い音声認識を実現することを目的とする。

【0006】

【課題を解決するための手段】

上記の目的を達成する本発明は、所定の音声と、予め収録された音声データの音韻隠れマルコフモデルとのマッチングを取ることにより音声認識を行う、次のように構成された音声認識装置として実現される。この音声認識装置は、認識対象である入力音声の特徴量を

10

20

30

40

50

抽出する特徴量抽出部と、予め収録された音声データの音韻隠れマルコフモデルと予め収録された雑音データの隠れマルコフモデルとを合成し合成モデルを作成する合成モデル作成部と、特徴量抽出部にて抽出された入力音声の特徴量と合成モデル作成部にて作成された合成モデルとのマッチングを取ることににより入力音声を認識する音声認識部とを備える。

【0007】

ここで、この音声認識部は、入力音声における発話区間を区切る適当な区間ごとに、具体的には例えば音声フレームごとに、独立してマッチング対象となる合成モデルを選択し、この入力音声の特徴量と選択された合成モデルとのマッチングを取ることを特徴とする。さらに、この音声認識装置において、音声認識部は、発話中の入力音声に含まれる雑音の

10

【0008】

また、本発明による他の音声認識装置は、音声認識のためのモデルとなる音声データを格納した音声データベースと、所定の雑音環境で発生が想定される雑音データを格納した雑音データベースと、音声データベースから読み出した音声データに基づき作成される音声モデルと雑音データベースから読み出した雑音データに基づき作成される雑音モデルとを合成して合成モデルを作成する合成モデル作成部と、認識対象である入力音声の特徴量と合成モデルとに関して、この入力音声の音声フレームごとに独立してマッチングを取ることににより音声認識を行う音声認識部とを備えることを特徴とする。

20

【0009】

さらにまた、本発明の他の音声認識装置は、音声認識のためのモデルとなる音声データを格納した音声データベースと、所定の雑音環境で発生が想定される雑音データを格納した雑音データベースと、音声データベースから読み出した音声データに基づき作成される音声モデルと雑音データベースから読み出した雑音データに基づき作成される雑音モデルとを合成して合成モデルを作成する合成モデル作成部と、認識対象である入力音声の特徴量と合成モデルとに関して、発話中の入力音声に含まれる雑音の変化に応じてマッチング対象となる合成モデルを動的に選択しながらマッチングを取ることににより音声認識を行う音声認識部とを備えることを特徴とする。

【0010】

また、上記の目的を達成する他の本発明は、コンピュータを制御して、音声を認識する、次のような音声認識方法として実現される。この音声認識方法は、認識対象である入力音声の特徴量を抽出し、メモリに格納するステップと、所定の音声データと雑音データとに基づいて生成された雑音の要素が加味された音韻隠れマルコフモデルをメモリから読み出すステップと、入力音声の音声フレームごとに入力音声の特徴量と音韻隠れマルコフモデルとのマッチングを取り、マッチングの結果に基づいて入力音声を認識するステップとを含むことを特徴とする。ここで、より詳しくは、入力音声を認識するステップは、入力音声の音声フレームごとに、独立してマッチング対象となる音韻隠れマルコフモデルを選択し、入力音声の特徴量と音韻隠れマルコフモデルとのマッチングを取る。

30

【0011】

さらに本発明による他の音声認識方法は、認識対象である入力音声の特徴量を抽出し、メモリに格納するステップと、所定の音声データと雑音データとに基づいて生成された雑音の要素が加味された音韻隠れマルコフモデルをメモリから読み出すステップと、入力音声の特徴量と音韻隠れマルコフモデルとに関して、発話中の入力音声に含まれる雑音の変化に応じてマッチング対象となる音韻隠れマルコフモデルを動的に選択しながらマッチングを取ることににより入力音声を認識するステップとを含むことを特徴とする。

40

【0012】

また、本発明は、コンピュータを制御して上記の音声認識装置を実現し、あるいは上記の音声認識方法の各ステップに対応する処理をコンピュータに実行させるプログラムとして実現される。このプログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記

50

録媒体に格納して配布したり、ネットワークを介して配信したりすることにより提供される。

#### 【0013】

##### 【発明の実施の形態】

以下、添付図面に示す実施の形態に基づいて、この発明を詳細に説明する。

図1は、本実施の形態による音声認識システムを実現するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

図1に示すコンピュータ装置は、演算手段であるCPU (Central Processing Unit: 中央処理装置) 101と、M/B (マザーボード) チップセット102及びCPUバスを介してCPU101に接続されたメインメモリ103と、同じくM/Bチップセット102及びAGP (Accelerated Graphics Port) を介してCPU101に接続されたビデオカード104と、PCI (Peripheral Component Interconnect) バスを介してM/Bチップセット102に接続されたハードディスク105及びネットワークインターフェイス106と、さらにこのPCIバスからブリッジ回路107及びISA (Industry Standard Architecture) バスなどの低速なバスを介してM/Bチップセット102に接続されたフロッピーディスクドライブ108及びキーボード/マウス109とを備える。また、処理対象である音声を入力し、音声データに変換してCPU101へ供給するためのサウンドカード (サウンドチップ) 110及びマイクロフォン111を備える。

なお、図1は本実施の形態を実現するコンピュータ装置のハードウェア構成を例示するに過ぎず、本実施の形態を適用可能であれば、他の種々の構成を取ることができる。例えば、ビデオカード104を設ける代わりに、ビデオメモリのみを搭載し、CPU101にてイメージデータを処理する構成としても良いし、ATA (AT Attachment) などのインターフェイスを介してCD-ROM (Compact Disc Read Only Memory) やDVD-ROM (Digital Versatile Disc Read Only Memory) のドライブを設けても良い。

#### 【0014】

図2は、図1に示したコンピュータ装置にて実現される本実施の形態による音声認識システムの構成を示す図である。

本実施の形態は、自動車の車室内のような発生する雑音の種類がある程度限定される環境において、HMM (隠れマルコフモデル) 合成法を用いて、突発的に発生する雑音や不規則に発生する雑音に対処し、高精度な音声認識を行う音声認識システムを実現する。

図2に示すように、本実施の形態による音声認識システムは、音声入力部10と、特徴量抽出部20と、音声認識部30とを備えると共に、音声認識部30にて使用される合成HMMを生成する合成HMM作成部40を備えている。また、合成HMM作成部40にて合成HMMを生成するために用いられる音声データベース50及び雑音データベース60を備える。

#### 【0015】

上記の構成において、特徴量抽出部20、音声認識部30及び合成HMM作成部40は、図1に示したメインメモリ103に展開されたプログラムにてCPU101を制御することにより実現される仮想的なソフトウェアブロックである。CPU101を制御してこれらの機能を実現させる当該プログラムは、磁気ディスクや光ディスク、半導体メモリ、その他の記録媒体に格納して配布したり、ネットワークを介して配信したりすることにより提供される。本実施の形態では、図1に示したネットワークインターフェイス106やフロッピーディスクドライブ108、図示しないCD-ROMドライブなどを介して当該プログラムを入力し、ハードディスク105に格納する。そして、ハードディスク105に格納されたプログラムをメインメモリ103に読み込んで展開し、CPU101にて実行することにより、図2に示した各構成要素の機能を実現する。

また、音声入力部10は、マイクロフォン111及びサウンドカード110にて実現される。音声データベース50及び雑音データベース60は、例えばハードディスク105にて実現される。

#### 【0016】

10

20

30

40

50

本実施の形態において、音声入力部 10 は、図 1 に示したマイクロフォン 111 及びサウンドカード 110 にて実現され、音声を入力すると共に、当該音声を電氣的な音声信号に変換して特徴量抽出部 20 に渡す。

特徴量抽出部 20 は、音声入力部 10 から受け取った音声信号に対して特徴量の抽出を行う。抽出された特徴量は、メインメモリ 103 や CPU 101 のキャッシュメモリにおける所定の領域に格納される。音声認識の HMM においては、音響パラメータとしてケプストラム係数が広く用いられており、特徴量抽出部 20 は、特徴量の抽出処理としてケプストラム分析を行うことができる。

音声認識部 30 は、特徴量抽出部 20 にて抽出された入力音声信号の特徴量と所定の音声モデル (HMM) とのマッチングを行い、マッチングの結果 (認識結果) として得られた文字 (テキスト) を出力する。本実施の形態では、後述する合成 HMM 作成部 40 にて作成される合成 HMM を用いることにより、音声認識部 30 による音声認識処理で使用される音響モデル (音韻モデル、単語モデルなど) を突発的な雑音や不規則な雑音の発生する環境に適應させてマッチングを行う。合成 HMM を用いたマッチングについては後述する。

#### 【0017】

合成 HMM 作成部 40 は、音声データベース 50 及び雑音データベース 60 にアクセスして合成 HMM を生成する。

図 3 は、合成 HMM 作成部 40 の機能を説明する機能ブロック図である。

図 3 を参照すると、合成 HMM 作成部 40 は、音声データベース 50 にアクセスして音声のモデル (音声 HMM) を作成する音声 HMM 作成部 41 と、雑音データベース 60 にアクセスして予め収録されている雑音のモデル (雑音 HMM) を作成する雑音 HMM 作成部 42 と、作成された音声 HMM と雑音 HMM とを合成して雑音の要素が加味された音韻隠れマルコフモデル (合成 HMM) を生成する HMM 合成部 43 とを備える。

#### 【0018】

音声データベース 50 には、雑音のない環境で収録された音声データが登録されており、音声 HMM 作成部 41 は、この音声データを用いて音声 HMM を作成する。作成された音声 HMM は、メインメモリ 103 や CPU 101 のキャッシュメモリの所定領域に保持される。

雑音データベース 60 には、本実施の形態における音声認識システムの使用環境で想定される雑音データが登録されており、雑音 HMM 作成部 42 は、この雑音データを用いて、雑音源ごとに独立に雑音 HMM を作成する。作成された雑音 HMM は、メインメモリ 103 や CPU 101 のキャッシュメモリの所定領域に保持される。

#### 【0019】

ここで、雑音データベース 60 について、さらに説明する。

実環境下では様々な雑音要因が存在するため、それら全てについて雑音データを収録し、雑音 HMM を作成しようとする、データ量が膨大になる。しかしながら、音声認識システムが使用される環境によっては、頻繁に発生する雑音の種類がある程度限定される場合がある。例えば、カーナビゲーションシステムの入力手段として搭載される音声認識システムの場合、車室内で頻繁に発生する雑音としては、比較的定常な走行中雑音 (エンジン音やロードノイズ) の他、非定常的な雑音としてマンホールなどを踏む音やウイカー音、ワイパーの動作する音などが想定される。そこで、音声認識システムが使用される環境に応じて、頻繁に発生することが想定される雑音について雑音データベース 60 を作成しておくことで、データ量が過大とならない実用的な音声認識システムを実現できる。なお、ハードウェア (CPU 101 等) の処理能力等に応じて、処理できる雑音データの量も変化することから、雑音データベース 60 のサイズを柔軟に変更できるのは言うまでもない。

#### 【0020】

HMM 合成部 43 は、音声 HMM 作成部 41 にて作成された音声 HMM と雑音 HMM 作成部 42 にて作成された雑音 HMM とをメインメモリ 103 等から取得し、これらを合成し

10

20

30

40

50

て合成HMMを作成する。

図4は、HMM合成部43の動作を説明する図である。

図4において、所定の音声/p/を構成するHMMの3つの状態のうち、i番目の状態の出力確率分布を $N_i(p)$ で表し、各雑音モデルの出力確率分布を $N(a)$ 、 $N(b)$ 、 $N(c)$ 、...で表す。

#### 【0021】

ここで、これらHMMにおける出力確率分布がケプストラム領域で作成されているものとする。この場合、HMM合成部43は、音声HMM及び雑音HMMのそれぞれに対し、コサイン変換を行い、さらに指数変換を行ってスペクトラル領域に変換した上で、分布の畳み込み(合成)を行う。

次に、合成された分布に対し、対数変換を行い、さらに逆コサイン変換を行ってケプストラム領域まで変換することにより、合成HMMを得る。分布の合成は、雑音源ごとに独立して用意された雑音HMMの出力確率分布に対して行われるため、合成HMMも雑音源ごとに独立に定義されることとなる。得られた合成HMMは、メインメモリ103やCPU101のキャッシュメモリにおける所定の領域に保持される。

以上の分布の変換については、例えば次の文献に詳細に記載されている。

文献:T. Takiguchi 他、"HMM-Separation-Based Speech Recognition for a Distant Moving Speaker," IEEE Transactions on speech and audio processing, Vol. 9, No. 2, pp. 127-140, 2001.

#### 【0022】

音声認識部30は、特徴量抽出部20にて抽出された入力音声信号の特徴量と、上記のようにして合成HMM作成部40により作成された合成HMMとのマッチングを取ることにより、この入力音声信号を認識する。なお、音声認識処理を完了するためには、本実施の形態にて行われる音響的な解析の他に、言語的な解析が行われることが必要であるが、この言語的な解析については本実施の形態による技術の対象ではなく、公知の技術を用いることができる。

ここで、本実施の形態における音声認識部30は、入力音声信号の特徴量と合成HMMとのマッチング(尤度計算)を、当該入力音声信号における音声フレーム単位で独立に行う。音声フレームとは、音声データにおける時間軸の最小単位である。

#### 【0023】

図5は、本実施の形態による音声認識部30の音声認識処理を説明するフローチャートである。

図5に示すように、音声認識部30は、メインメモリ103等から、特徴量抽出部20にて抽出された入力音声信号の特徴量と、上述した合成HMMとを取得し(ステップ501、502)、音声フレーム単位で、入力音声信号の特徴量との尤度が最大となる合成HMMを選択し(ステップ503)、その値をその時刻(音声フレーム)での尤度として採用する。採用された尤度は、メインメモリ103やCPU101のキャッシュメモリに一時的に保持される。

そして、発話終了まで、音声フレームごとに最も尤度が高くなる合成HMMを選択しながら、各時刻(音声フレーム)での最大尤度を加算していく(ステップ504)。すなわち、音声認識部30は、音声フレームについて最大尤度が得られたならば、メインメモリ103等に保持されている尤度を読み出して加算し、再びメインメモリ103等に保存する。これにより、ステップ503で選択された合成HMMの尤度が、直前の音声フレームまでの最大尤度の総和に随時加算されていく。この処理を発話終了まで繰り返すことにより、当該発話全体に対する尤度が算出される(ステップ505)。発話終了まで処理が尤度を加算する処理が行われたならば、算出された当該発話全体に対する尤度を用いて認識を行い、結果を出力する(ステップ505、506)。

#### 【0024】

以上のようにして、1つの発話に対する認識処理において、雑音を加味した合成HMMとのマッチングを音声フレーム単位で独立に行うことにより、突発的な雑音の発生などによ

10

20

30

40

50



り1つの発話中に雑音の状態や種類が変化した場合でも、マッチングにおいて適用する雑音モデルを動的に変更して対応することが可能となる。所定の入力音声信号において、どの部分が発話であるかについては、既存の手法を用いて判断することができる。

#### 【0025】

なお、上述した本実施の形態の動作においては、音声フレーム単位でマッチする（最大尤度の）合成HMMの探索を行ったが、一定の時間あるいは音声HMMの状態や音声HMMごとというように、発話区間を区切る適当な区間ごとに同一の雑音HMMを割り当てることにより、マッチングにおける合成HMMの探索時間を削減し処理コストを軽減することも可能である。この場合、非定常的な雑音に対する対応力は音声フレームごとにマッチングを行う場合に比べると低下するが、音声認識システムが使用される雑音環境（想定される雑音の種類等）に応じて適切な間隔を設定することにより、音声認識の精度を低下させることなく適用することができる。

#### 【0026】

次に、本実施の形態を用いた具体的な評価実験について説明する。

本実施の形態による音声認識システムを、自動車の車室内での音声認識に用い、雑音を考慮しない音声モデルを用いた認識（Clean HMMs）、従来のHMM合成法による認識（手法1）、本実施の形態による認識（手法2）で、認識率を測定し比較した。また、突発性の雑音としてハザード（ウィンカー）音、ある程度の時間長を持つ非定常雑音としてワイパーの動作音、定常雑音としてアイドリング時のエンジン雑音の3種類の雑音に対して本手法の有効性を検証した。

#### 【0027】

<評価1>

ここでは、評価音声データに、

- ・アイドリング時のエンジン雑音（以下、アイドリング雑音）
- ・ハザード（ウィンカー）音

の2種類の雑音を加算されている。ここで、ハザード音は、1周期が約0.4secである。

認識時に予め用意されている雑音HMMは、1.アイドリング雑音、2.走行雑音（一般道路を約40～50Kmで走行した際のロードノイズ）、3.ハザード音、4.ワイパー動作音、5.ハザード音+アイドリング雑音、6.ワイパー音+走行雑音の6種類である。また、1つの雑音HMMは、1つの状態と1つの多次元正規分布とで表されているものとする。

次に、これらの雑音HMMと音声HMM（55個の音韻HMM）との合成を行う。音声HMMは、各音韻が状態ごとに4つの多次元正規分布を持ち、この状態ごとに雑音HMMとの合成を行う。

信号の分析条件は、サンプリング周波数12kHz、フレーム幅32ms、分析周期8msである。音響特徴量としては、MFCC（Mel Frequency Cepstral Coefficient）16次元を用いた。また、テスト話者は男性1人で、500単語認識を行った。

#### 【0028】

図6は、以上の条件で行われた3種類の手法による音声認識の結果（認識率）を示す図表である。

図6を参照すると、定常的なアイドリング雑音しか対応できない手法1（従来のHMM合成法）に比べて、突発的なハザード音にも対応する手法2（本実施の形態）の方が、認識率が大きく改善されていることがわかる。

#### 【0029】

<評価2>

ここでは、評価音声データに、

- ・走行雑音（一般道路走行中）
- ・ワイパー動作音

の2種類の雑音を加算されている。ここで、ワイパー動作音は、1周期が約1.1sec

10

20

30

40

50

である。その他の条件は、＜評価１＞の条件と同じである。

図７は、以上の条件で行われた３種類の手法による音声認識の結果（認識率）を示す図表である。

図７を参照すると、上記の条件では、１つの発話中にワイパー動作音が発生している区間と無い区間とがあるため、発話区間内で適用する雑音モデルを動的に切り替える手法２（本実施の形態）の方が、手法１（従来のHMM合成法）よりも高い認識精度を得ていることがわかる。

【００３０】

【発明の効果】

以上説明したように、本発明によれば、突発的に発生する雑音や不規則に発生する雑音などのような急激な変化を伴う雑音に対しても十分に対処し、精度の高い音声認識を実現することができる。

10

【図面の簡単な説明】

【図１】 本実施の形態による音声認識システムを実現するのに好適なコンピュータ装置のハードウェア構成の例を模式的に示した図である。

【図２】 図１に示したコンピュータ装置にて実現される本実施の形態による音声認識システムの構成を示す図である。

【図３】 本実施の形態における合成HMM作成部の機能を説明する図である。

【図４】 本実施の形態におけるHMM合成部の動作を説明する図である。

【図５】 本実施の形態による音声認識部の音声認識処理を説明するフローチャートである。

20

【図６】 本実施の形態と従来の技術による音声認識の結果（認識率）を比較する図表である。

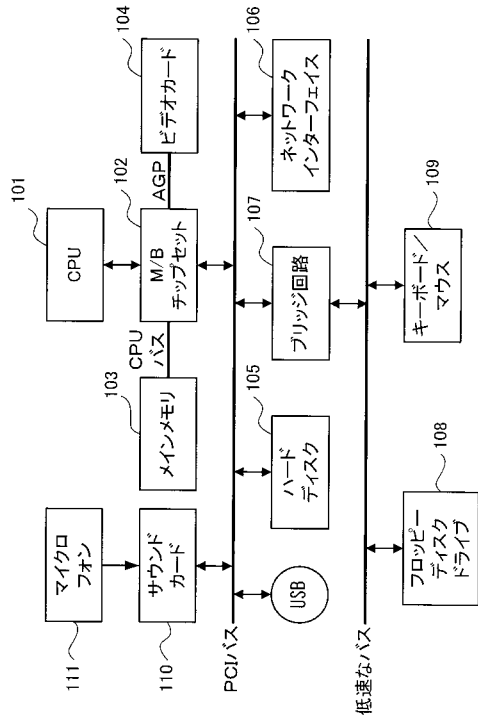
【図７】 本実施の形態と従来の技術による音声認識の他の結果（認識率）を比較する図表である。

【符号の説明】

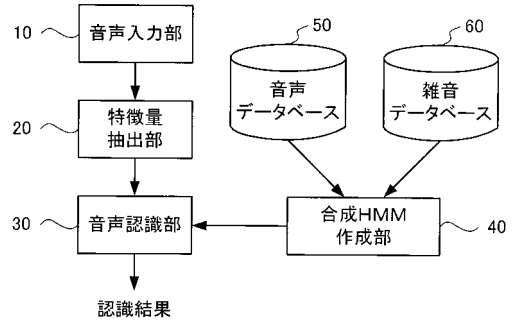
１０…音声入力部、２０…特徴量抽出部、３０…音声認識部、４０…合成HMM作成部、４１…音声HMM作成部、４２…雑音HMM作成部、４３…HMM合成部、１０１…CPU、１０２…M/Bチップセット、１０３…メインメモリ、１１０…サウンドカード、１１１…マイクロフォン

30

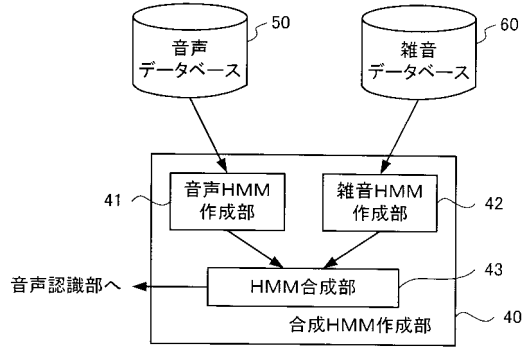
【図1】



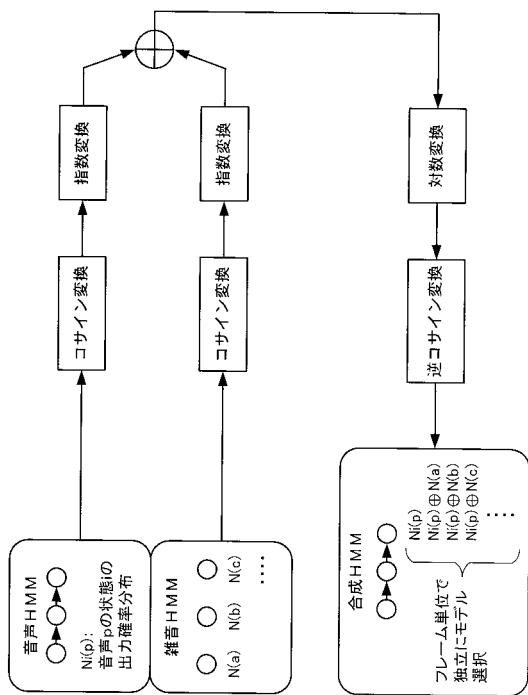
【図2】



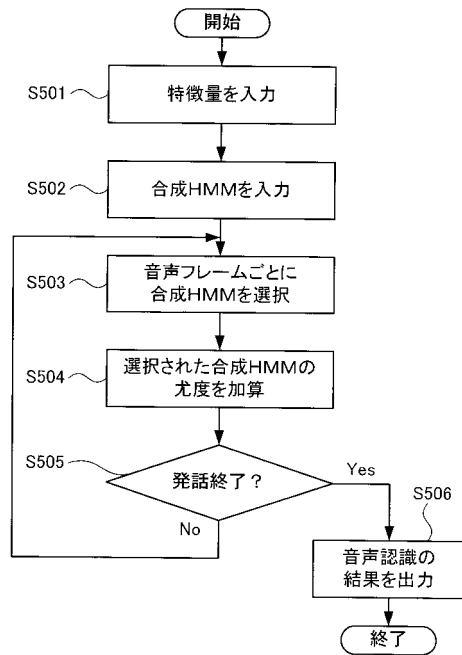
【図3】



【図4】



【図5】



## 【 図 6 】

Clean HMMs	手法 1	手法 2
61.2%	78.4%	90.2%

## 【 図 7 】

Clean HMMs	手法 1	手法 2
19.0%	61.4%	75.4%

---

フロントページの続き

- (72)発明者 滝口 哲也  
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内
- (72)発明者 西村 雅史  
神奈川県大和市下鶴間1623番地14 日本アイ・ピー・エム株式会社 東京基礎研究所内

合議体

- 審判長 原 光明  
審判官 加藤 恵一  
審判官 月野 洋一郎

- (56)参考文献 特開平10-11085(JP,A)  
特開2000-75889(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G10L15/00-15/28