

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4515054号
(P4515054)

(45) 発行日 平成22年7月28日 (2010. 7. 28)

(24) 登録日 平成22年5月21日 (2010. 5. 21)

(51) Int. Cl.	F I
G 1 0 L 15/02 (2006. 01)	G 1 0 L 15/02 3 0 0 F
G 1 0 L 15/10 (2006. 01)	G 1 0 L 15/10 3 0 0 G
G 1 0 L 15/18 (2006. 01)	G 1 0 L 15/18 2 0 0 E
	G 1 0 L 15/18 3 0 0 H

請求項の数 8 (全 26 頁)

(21) 出願番号	特願2003-278640 (P2003-278640)	(73) 特許権者	500046438
(22) 出願日	平成15年7月23日 (2003. 7. 23)		マイクロソフト コーポレーション
(65) 公開番号	特開2004-54298 (P2004-54298A)		アメリカ合衆国 ワシントン州 9805
(43) 公開日	平成16年2月19日 (2004. 2. 19)		2-6399 レッドモンド ワン マイ
審査請求日	平成18年7月19日 (2006. 7. 19)		クロソフト ウェイ
(31) 優先権主張番号	60/398, 166	(74) 代理人	100077481
(32) 優先日	平成14年7月23日 (2002. 7. 23)		弁理士 谷 義一
(33) 優先権主張国	米国 (US)	(74) 代理人	100088915
(31) 優先権主張番号	60/405, 971		弁理士 阿部 和夫
(32) 優先日	平成14年8月26日 (2002. 8. 26)	(72) 発明者	デン リ
(33) 優先権主張国	米国 (US)		アメリカ合衆国 98074 ワシントン
(31) 優先権主張番号	10/267, 522		州 サマミッシュ ノースイースト 30
(32) 優先日	平成14年10月9日 (2002. 10. 9)		ストリート 22310
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 音声認識の方法および音声信号を復号化する方法

(57) 【特許請求の範囲】

【請求項 1】

有限状態システムの現状態についてのスコアを生成することによって音声信号を復号化する方法であって、

先行する状態の終わりの最適な生成関係値に基づいて、前記現状態についての生成関係値を求めるステップであって、前記最適な生成関係値は、現状態についてのスコアを最大にする 1 組の連続的値内の生成関係値である、ステップと、

前記生成関係値を使用して音の尤度を判定するステップであって、該尤度は、前記先行する状態と前記現状態との間のパスに位置合せされる 1 組の観測ベクトルによって表されるステップと、

前記音の尤度を前記先行する状態からのスコアと組み合わせ、前記現状態についてのスコアを判定するステップであって、前記先行する状態からの前記スコアは軌跡の離散的クラスに関連し、該軌跡の離散的クラスにおいて可能な軌跡の連続値はクラスタ化され、該クラスは最適な生成関係値のクラスと合致する、ステップと

を備えたことを特徴とする音声信号を復号化する方法。

【請求項 2】

複数のクラス中の生成関係値の各クラスについて別々のスコアを生成するのに使用されることを特徴とする請求項 1 に記載の音声信号を復号化する方法。

【請求項 3】

前記生成関係値は、前記現状態での現在時刻と前記先行する状態の終わりとの間の時間

の長さにさらに基づくことを特徴とする請求項 1 又は 2 に記載の音声信号を復号化する方法。

【請求項 4】

前記生成関係値は、時間依存性補間重みを使用して計算されることを特徴とする請求項 1 から 3 のいずれかに記載の音声信号を復号化する方法。

【請求項 5】

前記時間依存性補間重みは、音に関連する時定数にさらに依存することを特徴とする請求項 4 に記載の音声信号を復号化する方法。

【請求項 6】

前記生成関係値を使用して音の尤度を判定するステップは、前記生成関係値に、音に依存しない値を掛けるステップを含むことを特徴とする請求項 5 に記載の音声信号を復号化する方法。

【請求項 7】

前記生成関係値を使用して音の尤度を判定するステップは、音がノイズであるとき、前記生成関係値に 0 を掛けるステップを含むことを特徴とする請求項 5 に記載の音声信号を復号化する方法。

【請求項 8】

音声信号 (u_j) の一部を表す観測値を受け取るステップと、
先行する音韻単位の終わりに対応する時間の軌跡を初期値 (g_0^j) とし、仮説音韻単位の生成関係ターゲット (T^{u_j}) との間の線形補間である生成関係ダイナミクス値を使用して、前記仮説音韻単位についての予測値

【数 1】

$$g_k^j = \beta^{u_j}(k) \cdot g_0^j + (1 - \beta^{u_j}(k)) \cdot T^{u_j}$$

を識別するステップであって、前記線形補間は時間依存性補間重み

【数 2】

$$\beta^{u_j}(k) = (1 + \gamma^{u_j} k) \cdot e^{-\gamma^{u_j} k}$$

を使用し、ここで、 u_j は前記音声信号 (u_j) に関連する時定数である、ステップと、
前記予測値 ($g(t)$) を特徴ベクトル ($o(t)$) にマッピングする音依存性マッピング関数 (h^{u_j}) を用いて、

【数 3】

$$z(t) = g(t) + w(t)$$

$$o(t) = h^{u_j}(z(t)) + v(t)$$

と表し、ここで、 $w(t)$ 、 $v(t)$ は前記観測値の期待値からの変位を表し、前記仮説音韻単位についての尤度

【数 4】

$$p(o(t)|z(t), u)$$

を判定するステップと

を備えたことを特徴とする音声認識の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はパターン認識に関連した音声認識の方法および音声信号を復号化する方法に関する。

【背景技術】

【0002】

音声認識システムなどのパターン認識システムは、入力信号を取得し、信号を復号化して、信号によって表されるパターンを見つけようと試みる。例えば音声認識システムでは、音声信号（しばしばテスト信号と呼ばれる）が認識システムによって受領され、復号化されて、音声信号によって表される単語の文字列が識別される。

10

【0003】

多くの音声認識システムでは、接続された状態の単一の列によって音声単位が表される隠れマルコフモデル（hidden Markov model）が使用される。トレーニング信号を使用して、状態を占有する確率分布と、状態間を遷移する確率分布が音声単位ごとに求められる。音声信号を復号化するために、信号がフレーム単位に分割され、各フレームは特徴ベクトルに変換される。次いで特徴ベクトルは各状態についての分布と比較され、フレームによって表すことのできる最も可能性の高い一続きのHMM（hidden Markov model）状態が識別される。次いでそのシーケンスに対応する音声単位が選択される。

20

【0004】

HMMベースの認識システムは、比較的単純な音声認識作業では良好に動作するが、音声のいくつかの重要な動的性質を直接にはモデル化しない（かつ、会話的な音声などの難しい作業についての動作は不十分であることが知られている）。その結果、HMMベースの認識システムは、トレーニングに使用される音声信号と、復号化する音声信号との動的な分節（articulation）の差を吸収することができない。例えば、くだけた発話設定では、話者は、自分の音声をハイポアーティキュレート（hypo-articulate）またはアンダーアーティキュレート（under-articulate）する傾向がある。これは、ユーザの音声分節の軌跡が、次のターゲットに再び向けられる前に、意図されるターゲットに達しない可能性があることを意味する。トレーニング信号は一般に、ハイポアーティキュレートされた音声よりも完全に分節された音声素材を話者が提供する、音声の「リーディング」スタイルを用いて形成されるので、ハイポアーティキュレートされた音声は、トレーニングされたHMM状態とは合致しない。その結果、レコグナイザ（recognizer）が提供する認識結果は、くだけた音声に対しては理想的な認識結果からほど遠い。

30

【0005】

同様の問題がハイパーアーティキュレートされた音声でも生じる。ハイパーアーティキュレートされた音声では、話者は、自分の音声の相異なる音を区別することができるように余分な労力を払う。この余分な労力には、似た音の音声単位とより良好に区別できるようにある音声単位の音を変更すること、ある音声単位の音をより長く保つこと、または各音が近隣の音と区別して知覚されるように、各音の間を急速に遷移させることが含まれる可能性がある。これらの各機構は、HMMシステムを使用して音声を認識することを難しくする。各技法により得られる、音声信号に対する1組の特徴ベクトルは、トレーニングデータ中に存在する特徴ベクトルとは十分に合致しないからである。

40

【0006】

HMMシステムはまた、発話する速度の変化に対処する上でも難点がある。したがって、誰かがトレーニング信号より遅く、または速く発話した場合、HMMシステムは、音声信号の復号化でより多くの誤りをする傾向がある。

【0007】

HMMシステムの代替方法が提案されている。具体的には、音声信号の生成関係パラメ

50

ータ (production-related parameter) の軌跡 (trajectory) または挙動を直接モデル化することが提案されている。しかし、諸提案のどれも、音声の動的性質を完全にはモデル化していない。具体的には、各モデルは、話者が音声単位に対する所望のターゲットに近づくときに生じる、時間に依存する軌跡の変化に対処していない。加えて、軌跡についての連続的値に基づく確率決定を可能にし、同時に管理可能な軌跡状態の数に検索スペースを制限する復号化手段を各モデルは提供していない。

【0008】

このことに照らして、音声軌跡の動的性質がより良好にモデル化され、復号化を管理可能とするように、他のモデル変数によって音声の生成関係挙動を明示的にモデル化する音声認識フレームワークが必要である。

10

【0009】

いくつかの文献に上述のような従来の技術に関連した技術内容が開示されている（例えば、非特許文献1参照）。

【0010】

【非特許文献1】L.Deng著「A Dynamic, Feature-Based Approach to the Interface Between Phonology and Phonetics for Speech Modeling and Recognition」Speech Communication、Vol. 24, No. 4, 1998年、pp.299-323

【発明の開示】

【発明が解決しようとする課題】

【0011】

20

従来のシステムには上述したような種々の問題があり、さらなる改善が望まれている。

【0012】

本発明は、このような状況に鑑みてなされたもので、その目的とするところは、容易に理想的な認識結果を得ることができる音声認識の方法および音声信号を復号化する方法を提供することにある。

【課題を解決するための手段】

【0013】

時間依存性補間重みを使用して、以前の時刻での生成関係ダイナミックス値とターゲットとの間の線形補間を実施することにより分節ダイナミックス値を識別する音声認識の方法が提供される。次いで生成関係ダイナミックス値を使用して予測音響特徴値を形成し、その予測音響特徴値を、観測した音響特徴値と比較して、観測した音響特徴値が所与の音韻単位によって生成された尤度を求める。

30

【0014】

ある実施形態では、以前の時刻での生成関係ダイナミックス値は、1組の連続的値から選択される。加えて、音韻単位の尤度が、以前の時刻での生成関係ダイナミックス値の離散的クラスに関連するスコアと組み合わせられ、現状態についてのスコアが求められる。

【発明の効果】

【0015】

以上説明したように本発明によれば、容易に理想的な認識結果を得ることができる。

【発明を実施するための最良の形態】

40

【0016】

以下、図面を参照して本発明の実施形態を詳細に説明する。

【0017】

図1に、本発明を実施することのできる適切なコンピューティングシステム環境100の一例を示す。コンピューティングシステム環境100は、適切なコンピューティング環境の一例に過ぎず、本発明の使用法または機能の範囲に関して何らかの制限を示唆するものではない。例示的動作環境100に図示する構成要素のうちのいずれか1つ、あるいはそれらの組合せに関係する何らかの依存関係または要件をコンピューティング環境100が有するものと解釈すべきでもない。

【0018】

50

本発明は、他の多数の汎用／専用コンピューティングシステム環境／構成で動作可能である。本発明と共に使用するのに適した周知のコンピューティングシステム、環境、および／または構成の例には、限定はしないが、パーソナルコンピュータ、サーバコンピュータ、ハンドヘルド／ラップトップ装置、マルチプロセッサシステム、マイクロプロセッサベースのシステム、セットトップボックス、プログラマブル家庭用電化製品、ネットワークPC、ミニコンピュータ、メインフレームコンピュータ、電話システム、ならびに上記のシステムまたは装置のいずれかを含む分散コンピューティング環境などが含まれる。

【0019】

本発明は、コンピュータが実行中の、プログラムモジュールなどのコンピュータ実行可能命令の一般的状況で説明することができる。一般に、プログラムモジュールは、特定のタスクを実行し、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、コンポーネント、データ構造などを含む。本発明はまた、通信ネットワークを介してリンクされるリモートプロセッサによってタスクが実行される分散コンピューティング環境でも実施することができる。分散コンピューティング環境では、プログラムモジュールは、メモリ記憶装置を含む、ローカルコンピュータ記憶媒体とリモートコンピュータ記憶媒体のどちらにも位置することができる。

【0020】

図1を参照すると、本発明を実施する例示的システムは、コンピュータ110の形態の汎用コンピューティング装置を含む。コンピュータ110の構成要素は、限定はしないが、プロセッサ120と、システムメモリ130と、システムメモリを含む様々なシステム構成要素をプロセッサ120に結合するシステムバス121とを含むことができる。システムバス121は、メモリバスまたはメモリコントローラと、周辺バスと、様々なバスアーキテクチャのうちのいずれかを用いるローカルバスとを含むいくつかのタイプのバス構造のうちのいずれでもよい。例えば、限定はしないが、このようなアーキテクチャには、ISA (Industry Standard Architecture) バス、MC A (Micro Channel Architecture) バス、EISA (Enhanced ISA) バス、VESA (Video Electronics Standards Association) ローカルバス、およびメザニンバスとも呼ばれるPCI (Peripheral Component Interconnect) バスが含まれる。

【0021】

コンピュータ110は、一般に様々なコンピュータ可読媒体を含む。コンピュータ可読媒体は、コンピュータ110がアクセス可能である入手可能などんな媒体でもよく、それには揮発性媒体と不揮発性媒体、リムーバブル媒体とノンリムーバブル媒体のどちらも含まれる。例えば、限定はしないが、コンピュータ可読媒体は、コンピュータ記憶媒体および通信媒体を含むことができる。コンピュータ記憶媒体には、コンピュータ可読命令、データ構造、プログラムモジュール、または他のデータなどの情報を格納するための何らかの方法または技術で実装される、揮発性媒体と不揮発性媒体、リムーバブル媒体とノンリムーバブル媒体のどちらも含まれる。コンピュータ記憶媒体には、限定はしないが、RAM、ROM、EEPROM (electrically erasable programmable read-only memory)、フラッシュメモリ、または他のメモリ技術、CD (compact disk) - ROM、デジタルバーサタイルディスク (DVD)、または他の光記憶装置、磁気力セット、磁気テープ、磁気ディスク記憶装置、または他の磁気記憶装置、あるいは、所望の情報を格納するのに使用することができ、コンピュータ110でアクセスすることができる他のどんな媒体も含まれる。通信媒体は一般に、コンピュータ可読命令、データ構造、プログラムモジュール、または他のデータを、搬送波または他の移送機構などの被変調データ信号でに実装し、その通信媒体にはどんな情報送達媒体も含まれる。「被変調データ信号」という用語は、信号の特性集合のうちの1つまたは複数を有する信号、または情報を符号化するように変化する信号を意味する。例えば、限定はしないが、通信媒体には、有線ネットワークまたは直接有線接続などの有線媒体、ならびに音響媒体、RF (radio frequencies) 媒体

10

20

30

40

50

、赤外線媒体、および他の無線媒体などの無線媒体が含まれる。上記のいずれの組合せも、コンピュータ可読媒体の範囲内に含まれるべきである。

【0022】

システムメモリ130は、読取り専用メモリ(ROM)131およびランダムアクセスメモリ(RAM)132などの揮発性メモリおよび/または不揮発性メモリの形態のコンピュータ記憶媒体を含む。起動中などにコンピュータ110内の要素間で情報を転送する助けになる基本ルーチンを含むBIOS(Basic Input/Output System)133が、一般にROM131内に格納される。一般に、プロセッサ120にとって即座にアクセス可能であり、かつ/またはプロセッサ120が現在操作しているデータおよび/またはプログラムモジュールをRAM132は含む。例えば、限定はしないが、図1にオペレーティングシステム134、アプリケーションプログラム135、他のプログラムモジュール136、およびプログラムデータ137を示す。

10

【0023】

コンピュータ110はまた、他の取外し可能/取外し不能な、揮発性/不揮発性コンピュータ記憶媒体も含むことができる。単なる一例であるが、図1に、ノンリムーバブル不揮発性磁気媒体を読み書きするハードディスクドライブ141と、リムーバブル不揮発性磁気ディスク152を読み書きする磁気ディスクドライブ151と、CD-ROMまたは他の光媒体などのリムーバブル不揮発性光ディスク156を読み書きする光ディスクドライブ155とを示す。例示的動作環境で使用するその他の取外し可能/取外し不能な揮発性/不揮発性コンピュータ記憶媒体には、限定はしないが、磁気テープカセット、フラッシュメモリカード、デジタルバーサタイルディスク、デジタルビデオテープ、ソリッドステートRAM、およびソリッドステートROMなどが含まれる。ハードディスクドライブ141は一般に、インタフェース140などのノンリムーバブルメモリインタフェースを介してシステムバス121に接続され、磁気ディスクドライブ151および光ディスクドライブ155は一般に、インタフェース150などのリムーバブルメモリインタフェースによってシステムバス121に接続される。

20

【0024】

上記で解説し、図1に図示するドライブとその関連するコンピュータ記憶媒体は、コンピュータ110に対してコンピュータ可読命令、データ構造、プログラムモジュール、および他のデータの記憶を実現する。例えば図1では、ハードディスクドライブ141がオペレーティングシステム144、アプリケーションプログラム145、他のプログラムモジュール146、およびプログラムデータ147を格納するものとして図示している。これらの構成要素は、オペレーティングシステム134、アプリケーションプログラム135、他のプログラムモジュール136、およびプログラムデータ137と同じであっても、異なってもよいことに留意されたい。オペレーティングシステム144、アプリケーションプログラム145、他のプログラムモジュール146、およびプログラムデータ147には、少なくともこれらが相異なるコピーであることを示すために異なる符号を付けてある。

30

【0025】

ユーザは、キーボード162と、マイクロフォン163と、マウス、トラックボール、またはタッチパッドなどのポインティングデバイス161などの入力装置を介して、コマンドおよび情報をコンピュータ110に入力することができる。他の入力装置(図示せず)には、ジョイスティック、ゲームパッド、サテライトディッシュ、スキャナなどを含めることができる。これらの入力装置や他の入力装置はしばしば、システムバスに結合されるユーザ入力インタフェース160を介してプロセッサ120に接続されるが、パラレルポート、ゲームポート、またはユニバーサルシリアルバス(USB)などの他のインタフェースおよびバス構造によって接続することもできる。モニタ191または他のタイプのディスプレイ装置もまた、ビデオインタフェース190などのインタフェースを介してシステムバス121に接続される。モニタに加えて、コンピュータは、スピーカ197やプリンタ196などの他の周辺出力装置も含むことができ、その周辺出力装置は、出力周辺

40

50

インタフェース 190 を介して接続することができる。

【0026】

コンピュータ 110 は、リモートコンピュータ 180 などの 1 つまたは複数のリモートコンピュータへの論理接続を使用して、ネットワーク環境で動作することができる。リモートコンピュータ 180 は、パーソナルコンピュータ、ハンドヘルド装置、サーバ、ルータ、ネットワーク PC、ピア装置、または他の共通ネットワークノードでよく、一般に、コンピュータ 110 に関して上記で述べた要素のうちの多数またはすべてを含む。図 1 に示す論理接続は、ローカルエリアネットワーク (LAN) 171 および広域ネットワーク (WAN) 173 を含むが、他のネットワークも含むことができる。このようなネットワーク環境は、オフィス、企業全体のコンピュータネットワーク、イントラネット、およびインターネットで一般的なものである。

10

【0027】

LAN ネットワーキング環境で使用する際、コンピュータ 110 は、ネットワークインタフェース / アダプタ 170 を介して LAN 171 に接続される。WAN ネットワーキング環境で使用する際、コンピュータ 110 は一般に、インターネットなどの WAN 173 を介して通信を確立するためのモデム 172 または他の手段を含む。モデム 172 は内蔵でも外付けでもよく、ユーザ入力インタフェース 160、または他の適切な機構を介してシステムバス 121 に接続することができる。ネットワーク環境では、コンピュータ 100 に関して示したプログラムモジュールまたはその一部を、リモートメモリ記憶装置内に格納することができる。例えば、限定はしないが、図 1 に、リモートアプリケーションプログラム 185 がリモートコンピュータ 180 上に常駐するものとして示す。図示するネットワーク接続は例示的なものであって、コンピュータ間の通信リンクを確立する他の手段も使用できることを理解されたい。

20

【0028】

図 2 は、例示的コンピューティング環境であるモバイル装置 200 のブロック図である。モバイル装置 200 は、マイクロプロセッサ 202、メモリ 204、入出力 (I/O) 構成要素 206、ならびにリモートコンピュータまたは他のモバイル装置と通信する通信インタフェース 208 を含む。一実施形態では、上記の構成要素は、適切なバス 210 を介して互いに通信するように結合される。

【0029】

30

メモリ 204 は、モバイル装置 200 への総電力が遮断されたときにメモリ 204 中に格納された情報が失われないようにバッテリバックアップモジュール (図示せず) を備えるランダムアクセスメモリ (RAM) などの不揮発性電子メモリとして実装される。メモリ 204 の一部は、プログラム実行用のアドレス指定可能メモリとして割り振ることが好ましく、メモリ 204 の別の部分は、ディスクドライブ上の記憶をシミュレートするための記憶用に使用することが好ましい。

【0030】

メモリ 204 は、オペレーティングシステム 212、アプリケーションプログラム 214、およびオブジェクトストア 216 を含む。動作中、オペレーティングシステム 212 をプロセッサ 202 によってメモリ 204 から実行することが好ましい。ある好適実施形態でのオペレーティングシステム 212 は、Microsoft Corporation から市販されている WINDOWS (登録商標) CE ブランドのオペレーティングシステムである。好ましくは、オペレーティングシステム 212 は、モバイル装置用に設計され、1 組の公開されたアプリケーションプログラミングインタフェース / メソッドを介してアプリケーション 214 が使用することのできるデータベース機能を実装する。オブジェクトストア 216 内のオブジェクトは、アプリケーション 214 およびオペレーティングシステム 212 によって維持され、公開されたアプリケーションプログラミングインタフェース / メソッドに対する呼出しに少なくとも部分的に応答する。

40

【0031】

通信インタフェース 208 は、モバイル装置 200 が情報を送受信することを可能にす

50

る多数の装置および技術を表す。この装置には、ほんの少数の例を挙げれば、有線／無線モデム、サテライト受信機、および放送チューナが含まれる。モバイル装置 200 はまた、コンピュータと直接接続して、それらの間でデータを交換することもできる。そのような場合、通信インタフェース 208 は、赤外線トランシーバまたはシリアル／パラレル通信接続でよい。これらのすべては、ストリーミング情報を伝送することができる。

【0032】

入出力構成要素 206 は、タッチセンシティブ画面、ボタン、ローラ、およびマイクロフォンなどの様々な入力装置と、オーディオジェネレータ、振動装置、およびディスプレイを含む様々な出力装置とを含む。上記で列挙した装置は例であり、かつモバイル装置 200 上にすべて存在する必要はない。加えて、他の入出力装置を、本発明の範囲内のモバイル装置 200 に取り付けることができ、または本発明の範囲内のモバイル装置 200 に関して見つけることができる。

【0033】

本発明は、音声の生成モデルを提供する。このモデルの下では、音声は、一続きの音韻単位の言語的定義を音声的に実施する話者による試行として表される。この試行の間、話者は、前の音韻単位の終わりに関連する値から、現在の音韻単位に関連するターゲットまでの軌跡をたどる生成関係値を生成する。この軌跡は、前の音韻単位の終わりの値と、ターゲット値の間の補間としてモデル化され、この補間は、2つの補間点間の重み付けが音韻単位の長さによって変化する時間依存性重みを含む。

【0034】

本発明のモデルは隠れ軌跡モデル (Hidden Trajectory Model) と呼ばれる。隠れ軌跡モデルは、ダイナミックスがリカーシブ (recursive) 形式によってではなく明示的形式の時間依存性によって規定されるという点で、隠れ動的モデル (Hidden-Dynamic Model) の特別な形である。この隠れ軌跡モデルは、隠れ生成関係パラメータ (ボーカルトラック共鳴など) を説明する動的／軌跡モデル構成要素と、生成関係パラメータを、メル周波数ケプストラル係数などの観測可能な音響特徴に変換するマッピングモデル構成要素の2つの層を含む。動的モデルは、対応する境界 ($1 \dots N+1$) と共に音シーケンス ($u_1, \dots, u_j, \dots, u_N$) が与えられた場合、生成関係パラメータについての一続きの軌跡値 ($z(1), \dots, z(t), \dots, z(T)$) を予測する。マッピングモデルは、その一続きの軌跡値、音シーケンス、および音境界が与えられた場合、一続きの音響観測ベクトルを予測する。

【0035】

どちらの層でも、統計モデルが与えられる。

【0036】

【数1】

$$z(t) = g(t) + w(t) \quad \text{式1}$$

$$o(t) = h^{u_j}(z(t)) + v(t) \quad \text{式2}$$

【0037】

上式で、 $g(t)$ は予測される軌跡であり、

【0038】

【数2】

$$h^{u_j}$$

【0039】

は生成関係パラメータを特徴空間にマッピングする音依存性マッピング関数である。

【0040】

加数 $w(t)$ および $v(t)$ は、平均値 0 を有し、実際の観測値の、期待値からの変位をモデル化する共分散行列 $Q = C_{ww}$ および $R = C_{vv}$ をそれぞれ有する i.i.d. ガウスノイズを表す。すなわち、

【 0 0 4 1 】

【 数 3 】

$$p(z(t)|g(t)) = N(z(t); g(t), Q) \quad \text{式 3}$$

$$p(o(t)|z(t), u) = N(o(t); h^{u_j}(z(t)), R) \quad \text{式 4}$$

【 0 0 4 2 】

本発明の下では、言語中の各音について、別々の 1 組の軌跡パラメータがトレーニングされる。加えて、音内の軌跡が、音の始まりでの軌跡の値と、音の開始から経過した時間とに基づいて決定される。したがって、ローカル時間 $k = t - \tau_j$ に対して、音 u_j 内の軌跡を説明することができる。ただし、 τ_j は音 u_j が始まる時間であり、各音は持続時間 $K^j = \tau_j + 1 - \tau_j$ を有する。

10

【 0 0 4 3 】

本発明の下では、任意のローカル時間インデックス k での、音 u_j についての軌跡は以下のように定義される。

【 0 0 4 4 】

【 数 4 】

$$g_k^j = \beta^{u_j}(k) \cdot g_0^j + (1 - \beta^{u_j}(k)) \cdot T^{u_j} \quad \text{式 5}$$

$$\beta^{u_j}(k) = (1 + \gamma^{u_j} k) \cdot e^{-\gamma^{u_j} k} \quad \text{式 6}$$

【 0 0 4 5 】

20

上式で、

【 0 0 4 6 】

【 数 5 】

$$T^{u_j}$$

【 0 0 4 7 】

は、音 u_j の軌跡に対するターゲットであり、

【 0 0 4 8 】

【 数 6 】

$$\gamma^{u_j}$$

30

【 0 0 4 9 】

は、音 u_j に関連する時定数であり、

【 0 0 5 0 】

【 数 7 】

$$g_0^j$$

【 0 0 5 1 】

は、その音に進入したときの軌跡の初期値である。式 6 の右辺の括弧で閉じられた項は、臨界ダンピング関数であることに留意されたい。本発明が、

40

【 0 0 5 2 】

【 数 8 】

$$g_0^j = g_{K^{j-1}}^{j-1} \quad \text{式 7}$$

【 0 0 5 3 】

を保証することによって音の間の連続性を実現することにも留意されたい。上式で、

【 0 0 5 4 】

【 数 9 】

$$g_0^j$$

50

【 0 0 5 5 】

は音 u_j の始めでの軌跡の値であり、

【 0 0 5 6 】

【 数 1 0 】

$g_{K^{j-1}}^{j-1}$

【 0 0 5 7 】

は音 u_{j-1} の終わりでの軌跡の値である。

【 0 0 5 8 】

式 5 および 7 は 2 つの重要な性質を有する。第 1 に、時刻 $k = 0$ では、軌跡は、先行する音についての軌跡の終わりの値と等しい。第 2 に、音が長期間にわたって持続する場合、軌跡はターゲット

10

【 0 0 5 9 】

【 数 1 1 】

T^{u_j}

【 0 0 6 0 】

に達する。したがって、式 5 で計算される軌跡は、音

【 0 0 6 1 】

【 数 1 2 】

20

g_0^j

【 0 0 6 2 】

の始まりでの生成関係ダイナミクス値と、生成関係ターゲット

【 0 0 6 3 】

【 数 1 3 】

T^{u_j}

【 0 0 6 4 】

との間の、時間依存性補間重み

30

【 0 0 6 5 】

【 数 1 4 】

$\beta^{u_j}(k)$

【 0 0 6 6 】

を用いた線形補間である。

【 0 0 6 7 】

本発明の下では、生成関係軌跡を音響特徴にマッピングするのに使用する音依存性マッピング関数

【 0 0 6 8 】

40

【 数 1 5 】

h^{u_j}

【 0 0 6 9 】

は、以下のように定義される区分的線形関数である。

【 0 0 7 0 】

【 数 1 6 】

$h^{u_j}(z(t)) = H_m Z(t) + h_m$

式 8

【 0 0 7 1 】

50

上式で m は、フレームまたは音全体に対して一定である混合インデックスである。一実施形態では、 H_m および h_m は音に依存しない。別の実施形態では、 H_m および h_m は、音、音クラス、または左から右に配置されたHMM状態に対応するサブフォン(sub-phone)ユニットに依存するようにされる。

【0072】

本発明の一態様の下では、静寂およびノイズに対して予測されるベクトルが、 $H_m = 0$ と仮定することによって形成される。その結果、静寂およびノイズに対して予測される特徴ベクトルは、生成関係値の軌跡に依存しない。これは、静寂およびノイズが音声生成での中断を表す生成モデルに適合する。

【0073】

このマッピング関数を使用すると、式2および4は以下のようになる。

【0074】

【数17】

$$o(t) = H_m z(t) + h_m + v(t) \quad \text{式9}$$

$$p(o(t)|z(t), u) = N(o(t); H_m z(t) + h_m, R_m) \quad \text{式10}$$

【0075】

モデルパラメータ

【0076】

【数18】

$$T^{u_i}, \gamma^{u_i}$$

【0077】

H_m 、 h_m 、 Q 、 R_m は、予測最大化トレーニングアルゴリズムを使用してトレーニングされる。このアルゴリズムはEステップを含む。Eステップでは、1組のトレーニング観測ベクトルがモデルパラメータの初期推定と共に使用され、十分な統計が発生し、混合重み、軌跡、および軌跡の2乗を含むある隠れ変数の値が予測される。具体的には、Eステップは以下の計算を含む。

【0078】

【数19】

$$\omega_m^n = P(m|o^n) = \frac{P(o^n|m)P(m)}{\sum_{m'=1}^M P(o^n|m')P(m')} \quad \text{式11}$$

$$p(o^n|m) = \prod_{k=1}^{K^n} p(o_k^n|m) = \prod_{k=1}^{K^n} N(o_k^n; H_m g_k^n + h_m, S_m) \quad \text{式12}$$

$$E[z_k^n | o^n, m] = [H_m^{\text{TRANS}} R_m^{-1} H_m + Q^{-1}]^{-1} [H_m^{\text{TRANS}} R_m^{-1} (o_k^n - h_m) + Q^{-1} g_k^n] \quad \text{式13}$$

$$\begin{aligned} E[z_k^n z_k^{n \text{trans}} | o^n, m] &= E[z_k^n z_k^{n \text{trans}} | o_k^n, m] \\ &= [H_m^{\text{TRANS}} R_m^{-1} H_m + Q^{-1}]^{-1} + E[z_k^n | o^n, m] E[z_k^{n \text{trans}} | o_k^n, m]^{\text{Trans}} \quad \text{式14} \end{aligned}$$

【0079】

上式で、

【0080】

【数20】

$$S_m = H_m Q H_m^{\text{trans}} + R_m \quad \text{式15}$$

【0081】

【数 2 1】

Θ_m^n

【0 0 8 2】

は、音の境界によって定義される、トークン n についての混合重みであり、

【0 0 8 3】

【数 2 2】

O_k^n

【0 0 8 4】

は、トークン n についての k 番目に観測されるトレーニングベクトルであり、

【0 0 8 5】

【数 2 3】

g_k^n

【0 0 8 6】

は、ローカル時間点 k でのトークン n 中の予測される軌跡の値であり、

【0 0 8 7】

【数 2 4】

Z_k^n

【0 0 8 8】

は、ローカル時間点 k でのトークン n 中の実際の軌跡の値であり、

m はトークン n と共に使用する混合成分であり、

M は、トークン n に関連する混合成分の数であり、

各混合成分の確率 $P(m)$ は一様で、 $1/M$ に等しく、

「trans」は行列の転置を表し、

【0 0 8 9】

【数 2 5】

$E[x|y]$

【0 0 9 0】

は、 y が与えられた場合の x の期待値を表す。

【0 0 9 1】

上記は EM アルゴリズム中の E ステップを完了する。M ステップの最初の反復を実施するために、モデルパラメータの初期推定を与えなければならない。一実施形態の下では、

【0 0 9 2】

【数 2 6】

T^{u_j}

【0 0 9 3】

および

【0 0 9 4】

【数 2 7】

$\gamma_i^{u_j}$

【0 0 9 5】

についての初期推定は、Klatt 音声合成器と、あるスペクトル分析結果とを組み合わせた情報を用いて選択される。次いで、周知の位置合せ用の既存の HMM システム / 技法を用いて、1 組のトレーニング観測ベクトルについて位置合せ境界 ($\gamma_1 \dots \gamma_{N+1}$)

10

20

30

40

50

）が決定される。ターゲット

【 0 0 9 6 】

【 数 2 8 】

T^{u_j}

【 0 0 9 7 】

、時定数

【 0 0 9 8 】

【 数 2 9 】

γ^{u_j}

10

【 0 0 9 9 】

、および位置合せ境界 (\dots) を使用して、1組の軌跡 $g(t)$ が上記の式5を使用して推定される。当初、各軌跡を求める際にノイズが0であり、かつ観測ベクトルを混合に対してランダムに割り当てると仮定すると、 H_m および h_m は、観測される特徴ベクトル $o(t)$ と計算される特徴ベクトルとの間の平方誤差の和を最小にするように推定される。ただし誤差は以下のように計算される。

【 0 1 0 0 】

$$v(t) = o(t) - (H_m g(t) + h_m) \quad \text{式 1 6}$$

各混合について $g(t)$ 、 H_m 、および h_m を決定した後、各混合について共分散行列 R_m を以下のように推定することができる。 20

【 0 1 0 1 】

【 数 3 0 】

$$\begin{aligned} R_m &= E\{vv^{\text{trans}}\} \\ &= \frac{1}{T} \sum_t v(t) \cdot v(t)^{\text{trans}} \\ &= \frac{1}{T} \sum_t (o(t) - (H_m g(t) + h_m)) \cdot (o(t) - (H_m g(t) + h_m))^{\text{trans}} \end{aligned} \quad \text{式 1 7}$$

【 0 1 0 2 】

30

Q の推定は、観測ノイズ $v(t)$ が最小となるように、軌跡 $g(t)$ の決定の際にノイズ $w(t)$ をまず推定することによって決定される。これにより、以下の式が得られる。

【 0 1 0 3 】

【 数 3 1 】

$$w(t) = (H_m^{\text{trans}} H_m)^{-1} H_m^{\text{trans}} \cdot (o(t) - (H_m g(t) + h_m)) \quad \text{式 1 8}$$

【 0 1 0 4 】

次いで Q を以下のように初期化する。

【 0 1 0 5 】

【 数 3 2 】

40

$$\begin{aligned} Q &= E\{ww^{\text{trans}}\} \\ &= \frac{1}{T} \sum_t w(t) \cdot w(t)^{\text{trans}} \end{aligned} \quad \text{式 1 9}$$

【 0 1 0 6 】

一実施形態の下では、 R_m および Q は対角行列であると仮定され、したがって行列の対角成分だけが計算される。

【 0 1 0 7 】

初期モデルパラメータを上記のように得た後、それらをEステップの結果と共に使用して、Mステップでのモデルパラメータを推定する。具体的には、モデルパラメータは以下 50

のように計算される。

【 0 1 0 8 】

【 数 3 3 】

$$\hat{T} = \frac{\sum_{n=1}^N \sum_{k=1}^{K^n} \sum_{m=1}^{M^n} \omega_m^n \left(E \left[z_k^n \middle| o^n, m \right] - g_o^n(k) \right)}{\sum_{n=1}^N K^n \cdot \sum_{m=1}^M \omega_m^n}$$

式 2 0

【 0 1 0 9 】

上式で、

【 0 1 1 0 】

【 数 3 4 】

\hat{T}

【 0 1 1 1 】

は T、

【 0 1 1 2 】

【 数 3 5 】

$$g_o^n(k) = (g_{K^{j-1}} - T)(1 + \gamma k) \exp(-\gamma k)$$

10

20

【 0 1 1 3 】

に対する更新であり、

および T は、以前に推定したターゲットである。

【 0 1 1 4 】

【 数 3 6 】

$$\hat{H}_m = (A - CB) \cdot (I - DB)^{-1}$$

式 2 1

【 0 1 1 5 】

ただし、

【 0 1 1 6 】

【 数 3 7 】

$$A = \left\{ \sum_{n=1}^N \omega_m^n \sum_{k=1}^{K^n} \left(o_k^n \cdot E \left\{ z_k^n \right\}^{\text{trans}} \right) \right\} \cdot X$$

式 2 2

$$B = \left\{ \sum_{n=1}^N \omega_m^n \sum_{k=1}^{K^n} E \left\{ z_k^n \right\}^{\text{trans}} \right\} \cdot X$$

式 2 3

$$C = \left\{ \sum_{n=1}^N \omega_m^n \sum_{k=1}^{K^n} \left(o_k^n \right) \right\} \cdot Y$$

式 2 4

$$D = \left\{ \sum_{n=1}^N \omega_m^n \sum_{k=1}^{K^n} E \left\{ z_k^n \right\} \right\} \cdot Y$$

式 2 5

$$X = \left\{ \sum_{n=1}^N \omega_m^n \sum_{k=1}^{K^n} E \left\{ z_k^n z_k^{n \text{trans}} \right\} \right\}^{-1}$$

式 2 6

$$Y = \frac{1}{\sum_{n=1}^N K^n \omega_m^n}$$

式 2 7

【 0 1 1 7 】

であり、I は恒等行列である。

30

40

50

【 0 1 1 8 】

【 数 3 8 】

$$\hat{h}_m = \frac{\sum_{n=1}^N \left[\omega_m^n \sum_{k=1}^{K^n} \left(o_k^n - \hat{H}_m E[z_k^n | o^n, m] \right) \right]}{\sum_{n=1}^N K^n \omega_m^n} \quad \text{式 2 8}$$

$$\hat{R}_m = \frac{\sum_{n=1}^N \left[\omega_m^n \sum_{k=1}^{K^n} \left(o_k^n - \hat{H}_m E[z_k^n | o^n, m] - \hat{h}_m \right) \left(o_k^n - \hat{H}_m E[z_k^n | o^n, m] - \hat{h}_m \right)^{\text{trans}} \right]}{\sum_{n=1}^N K^n \omega_m^n} \quad \text{式 2 9}$$

$$\hat{Q} = \frac{\sum_{n=1}^N \left[\sum_{k=1}^{K^n} \left(E[z_k^n | o^n] - g_k^n \right) \left(E[z_k^n | o^n] - g_k^n \right)^{\text{trans}} \right]}{\sum_{n=1}^N K^n} \quad \text{式 3 0} \quad 10$$

【 0 1 1 9 】

かつ

【 0 1 2 0 】

【 数 3 9 】

 $\hat{\gamma}$

【 0 1 2 1 】

20

は勾配降下アルゴリズムを使用して求められる。ただし

【 0 1 2 2 】

【 数 4 0 】

 $\hat{\gamma}$

【 0 1 2 3 】

は、その変化が反復の間のしきい値よりも小さくなるまで漸進的に更新される。具体的には、

【 0 1 2 4 】

【 数 4 1 】

30

 $\hat{\gamma}$

【 0 1 2 5 】

は

【 0 1 2 6 】

【 数 4 2 】

$$\hat{\gamma}^{r+1} = \hat{\gamma}^r + \varepsilon (-2) Q^{-1} \sum_{n=1}^N \sum_{k=1}^{K^n} (z_k^n - g^n) \cdot (\hat{\gamma}^r \cdot k^2 \cdot (g_0^n - \hat{T}) \cdot e^{-ik}) \quad \text{式 3 1}$$

【 0 1 2 7 】

40

を使用して更新される。

【 0 1 2 8 】

EステップおよびMステップは、最終の1組のモデルパラメータに達するように何回か反復することができる。その最終の1組のパラメータを求めた後、それを使用して、観測した1組の音響ベクトルを復号化することができる。

【 0 1 2 9 】

復号化作業は、一続きの音響観測を生成した可能性の最も高い音のシーケンスを見つけるものである。本発明の一実施形態の下では、2つの状態を接続するエッジによって各音が表される有限状態変換器(transducer)を使用して復号化が実施される。ただし、時間内の状態の位置は、復号化中に決定される。したがって、復号化プロセスは、エッジのシ

50

ーケンス ($E = (e_1, \dots, e_j, \dots, e_n)$) と、音響観測値のシーケンス ($O = (o_1, \dots, o_j, \dots, o_L)$) を生成した可能性が最も高い状態位置とを見つけるものである。

【0130】

【数43】

$$\hat{E} = \arg \max_E P(E|O) \quad \text{式32}$$

【0131】

上式は、ベイズの公式を使用し、分母を無視すると、以下のようになる。

【0132】

【数44】

$$\hat{E} = \arg \max_E P(O|E)P(E) \quad \text{式33}$$

【0133】

式33の確率は以下のように求められる。

【0134】

【数45】

$$\arg \max_E P(O|E)P(E) = \max_{n \in \{n_{\text{term}}\}} \max_g \{H_{c(g)}(L, n)\} \quad \text{式34}$$

【0135】

上式で、

n_{term} は、すべての終端状態の集合であり、

L は、最後の観測の時間インデックスであり、

$H_{c(g)}(L, n)$ は、軌跡 $c(g)$ のクラス内の軌跡 g を有する、時刻 L でのノード n への最高のスコアリングパスである。ただし、最高のスコアリングパスは以下のように求められる。

【0136】

【数46】

$$H_{c(g)}(T, n) = \max_{e: n_{e2}=n} \max_{\Delta t} \{Q_{c(g), \Delta t}(t, e)\} \quad \text{式35}$$

【0137】

上式で、

$e: n_{e2}=n$ は、状態 n で終了するエッジの集合であり、

t は $t - t'$ に等しい。ただし t' は、現在の音またはエッジが始まった時間インデックスであり、 t は状態 n についての時間インデックスであり、したがって t は、現在の音についての可能な持続時間であり、

【0138】

【数47】

$$Q_{c(g), \Delta t}(t, e) = \max_{g'} \{\delta(c(g) = c(G_{g'}(\Delta t, e))) \cdot R_{c(g), \Delta t}(t, e)\} \quad \text{式36}$$

【0139】

上式で、

$Q_{c(g)}, t(t, e)$ は、時刻 $t' = t - t$ にエッジに進入し、エッジ e に沿うクラス $c(g)$ の軌跡に沿って時刻 t にノード n に進入する最良のパスの確率と理解することができる。

【0140】

g' は、前のエッジの終わりまでの軌跡の履歴であり、

$c(\cdot)$ は、可能な軌跡の連続的な値がクラスタ化された、軌跡の離散的クラスを示し、

【0141】

10

20

30

40

50

【数 4 8】

$$\delta(A=B)_E = \begin{cases} 1 & A=B \text{ の場合} \\ 0 & \text{上記以外の場合} \end{cases} \quad \text{式 3 7}$$

$$G_{g'}(\Delta t, e) = \beta(\Delta t) \cdot g' + (1 - \beta(\Delta t)) \cdot T \quad \text{式 3 8}$$

【 0 1 4 2 】

【数 4 9】

$$R_{c(g'), \Delta t}(t, e) = \left\{ p(o(t) | \Delta t, e, g'_{\text{opt}}) \frac{P(\Delta t | e)}{P(\Delta t - 1 | e)} \cdot R_{c(g'), \Delta t - 1}(t - 1, e) \right\} \quad \text{式 3 9} \quad 10$$

【 0 1 4 3 】

ただし

【 0 1 4 4 】

【数 5 0】

$$R_{c(g'), 0}(t, e) = P(\Delta t = 0 | e) \cdot P(e) \cdot H_{c(g')}(t, n_{\text{el}}) \quad \text{式 4 0}$$

$$p(o(t) | \Delta t, e, g'_{\text{opt}}) = \sum_m P(m | u, g) \cdot N(o(t); H_m g_{\Delta t}^e + h_m; S_m) \quad \text{式 4 1}$$

【 0 1 4 5 】

20

これは、エッジ e に関連する仮説音についての尤度であり、仮説持続期間は以下のように近似される。

【 0 1 4 6 】

【数 5 1】

$$p(o(t) | \Delta t, e, g'_{\text{opt}}) \approx C \max_m \{ N(o(t); H_m g_{\Delta t}^e + h_m; S_m) \} \quad \text{式 4 2}$$

【 0 1 4 7 】

上式で、

S_m は上記の式 1 5 で定義される。

【 0 1 4 8 】

30

【数 5 2】

 g'_{opt}

【 0 1 4 9 】

は、

【 0 1 5 0 】

【数 5 3】

 $g_{\Delta t}^e$

【 0 1 5 1 】

40

を計算するのに使用するとき、式 3 6 の右辺を最大にする g' であり、現在の音についての開始生成関係値を表す。

【 0 1 5 2 】

【数 5 4】

 $g_{\Delta t}^e$

【 0 1 5 3 】

は、 t 、

【 0 1 5 4 】

【数 5 5】

g_{opt}^*

【 0 1 5 5】

、およびエッジ e に関連する軌跡パラメータが与えられたときの、生成関係ダイナミックス値の期待値であり、

C は定数であり、

【 0 1 5 6】

【数 5 6】

$H_m g_{\Delta t}^e + h_m$

10

【 0 1 5 7】

は、観測ベクトルについての予測値を表す。ただし H_m および h_m はエッジ e に関連する音に依存しない。

【 0 1 5 8】

式 4 2 の右辺の計算は、観測ベクトルと

【 0 1 5 9】

【数 5 7】

$H_m g_{\Delta t}^e + h_m$

20

【 0 1 6 0】

の差を取ることを含む。

【 0 1 6 1】

式 3 9 は、エッジ ($t' = t - t$) に割り当てられた最初の時間インデックスから現在の時間インデックス t までの各時間インデックスについて、単一の音を表すエッジ e 内で反復的に実施される。式 4 0 に示すように、最初の時間インデックスに達したとき、現在のエッジの最初の状態と、現在のエッジの最初の時間インデックスとを使用して式 3 5 を評価することによって別の反復が実施される。この反復は、式 4 0 の最後の項によって示される。

【 0 1 6 2】

30

式 3 9 の持続時間モデル、

【 0 1 6 3】

【数 5 8】

$P(\Delta t|e)$

【 0 1 6 4】

、および

【 0 1 6 5】

【数 5 9】

$P(\Delta t - 1|e)$

40

【 0 1 6 6】

は、トレーニング観測シーケンスをセグメント化する従来技術の HMM モデルを使用してトレーニングすることができる。セグメント化したテキストを使用して、各音についての持続期間が識別され、その持続期間は、各音についての持続期間確率分布を構築するのに使用される。

【 0 1 6 7】

上記のリカーシブなフレームワークで示したが、図 3 のフローチャートを参照しながら以下で説明するように、式 3 2 ~ 4 2 の復号化は、本発明の一実施形態の下では反復的に実施される。図 3 に示す方法では、観測される信号の各フレームで、1 組の最良のパスス

50

コアが求められる。具体的には、特定のフレームで利用可能な各状態について、上記の式 35 で求められる別々の $H_{c(g)}(t, n)$ 値が、軌跡の各クラスについて求められる。例えば 4 つの状態 400、402、404、および 406 を含む図 4 の単純な状態図の場合、軌跡のクラスが 3 つだけ存在することを条件として、図 5 の $H_{c(g)}(t, n)$ 値が生成される。したがって、値 500、502、および 504 は、第 1 フレーム（時間インデックス t_1 ）の間に、それぞれ状態 n_1 のクラス $c(g)1$ 、 $c(g)2$ 、および $c(g)3$ について生成される。値 506、508、および 510 は、第 2 フレーム（時間インデックス t_2 ）の間に、それぞれ状態 n_1 のクラス $c(g)1$ 、 $c(g)2$ 、および $c(g)3$ について生成される。図 6 に、観測されるベクトルの第 4 フレームについてのスコアを計算中の図 5 のダイアグラムを示す。図 6 は、以下の図 3 の説明で参照する。

10

【0168】

図 3 のステップ 300 では、最初に観測したベクトルを選択し、時間インデックスを 1 に設定する。ステップ 302 では、時間インデックスの識別中に利用可能な状態を識別する。状態は、時間インデックスがその状態と有限状態システムの開始状態 n_0 の間のエッジ数に等しい場合、その時間インデックスで利用可能である。例えば、図 4 の単純化した状態システムでは、4 つの状態 400、402、404、および 406 が存在する。時間インデックス 1 の間、利用可能な状態は 402 および 404 だけである。状態 400 と状態 404 の間の各エッジについて少なくとも 1 つの特徴ベクトルを有するための十分な観測された特徴ベクトルが存在しないからである。時間インデックス 2 の間、状態 402、404、および 406 はすべて利用可能となる。

20

【0169】

他の実施形態では、時間インデックスが、状態と開始状態 n_0 の間のエッジ数よりも大きいある量である場合、その状態は利用可能とみなされない。これにより、生じる可能性の低い状態についての一時的な位置をなくすることによって実行される計算の数が削減される。

【0170】

利用可能な状態を識別した後、ステップ 304 でその状態のうちの 1 つを選択する。次いでステップ 306 で、軌跡のクラス $c(g)$ を選択する。したがって図 6 では、ステップ 304 および 306 により、状態 n_2 およびクラス $c(g)1$ が選択される。これは、 $H_{c(g)}(t, n)$ 値を生成する対象の状態およびクラスである。

30

【0171】

ステップ 308 では、選択した状態に向かうエッジを選択する。例えば、状態 406 に対して図 4 のエッジ e_3 を選択する。エッジを選択した後、 t をステップ 310 で選択する。この選択により、特定の時間インデックスでの以前の状態が識別される。したがって例えば、ステップ 308 および 310 でエッジ e_3 および $t = 3$ を選択することにより、図 6 の前状態のような、時間インデックス t_1 での状態 n_1 の発生が選択される。

【0172】

ステップ 312 では、現在選択している軌跡のクラス、選択したエッジ、および選択した t について、式 36 に示す計算を最大にする g' を識別する。式 39 および 40 に示すように、式 36 の計算は、 t が交差する時間インデックスのすべてに沿ってリカーシブに実行され、選択したエッジの始めの状態に関連する $H_{c(g')}(t - t, n_{e_1})$ の値に依存する。値 $H_{c(g')}(t - t, n_{e_1})$ は、以前の時間インデックスで図 3 の方法を使用して先に計算されていること、および値は、式 36 にとって最適な g' のクラスに基づいて選択されることに留意されたい。したがって、 g' がクラス $c(g)3$ 内にある場合、図 6 の直線 600 で示すようにエッジ e_3 および $t = 3$ を使用して、 $H_{c(g)3}(t_1, n_1)$ が、時刻 t_4 での状態 N_2 についての式 39 および 40 の計算で使用される。

40

【0173】

ステップ 314 では、この方法は、考慮すべき追加の t が存在するかどうかを判定する。考慮すべき追加の t が存在する場合、プロセスはステップ 310 に戻り、異なる

50

t を選択する。次いで、新しい t を使用してステップ 312 を反復し、式 36 を使用して新しいスコアを識別する。

【0174】

t がもはや存在しないとき、式 36 を使用して計算したスコアを互いに比較し、ステップ 316 で、最高のスコアを現在のエッジについてのスコアとして選択する。例えば、図 6 の直線 600、602、および 604 で表されているように、エッジ e_3 に関して、 $t = 1$ 、 $t = 2$ 、 $t = 3$ のそれぞれについて別々のスコアを計算する。次いで最高のスコアをステップ 316 で選択し、それによって状態と、観測した音声フレームの間の特定の位置合せを選択する。

【0175】

ステップ 318 では、この方法は、現状態で終了するさらに別のエッジが存在するかどうかを判定する。さらに別のエッジが存在する場合、プロセスはステップ 308 に戻り、ステップ 310、312、314、および 316 を反復して、新しく選択したエッジについてのエッジスコアを生成する。例えば、図 6 では、エッジ e_4 についてエッジスコアが計算される。

【0176】

ステップ 318 で、現状態で終了するエッジがもはや存在しない場合、最高のエッジスコアを、選択したクラスおよび状態に対する $H_c(g)(t, n)$ パススコアとして選択する。

【0177】

ステップ 322 では、プロセスは、 $H_c(g)(t, n)$ パススコアを必要とする追加のクラスが存在するかどうかを判定する。追加のクラスが存在する場合、プロセスはステップ 306 に戻り、次の軌跡のクラスを選択する。次いでステップ 308 ~ 320 を反復し、新しく選択したクラスについて $H_c(g)(t, n)$ パススコアを生成する。ステップ 322 で別のクラスが存在しない場合、プロセスは、スコアリングする必要のある別の利用可能な状態が存在するかどうかをステップ 324 で判定する。別の利用可能な状態が存在する場合、プロセスはステップ 304 に戻って次の利用可能な状態を選択し、ステップ 306 ~ 322 を反復して、その状態についての $H_c(g)(t, n)$ パススコアの集合を生成する。

【0178】

現在の時間インデックスで利用可能な状態のすべてについてスコアを生成したとき、プロセスは、別の観測ベクトルが存在するかどうかをステップ 326 で判定する。別の観測ベクトルが存在する場合、プロセスはステップ 300 に戻り、ステップ 300 で次の観測ベクトルを選択し、時間インデックスを 1 つ増分する。次いで、新しいベクトルおよび時間インデックスについてステップ 302 ~ 324 を反復する。最後の観測ベクトルをステップ 326 で処理したとき、クラスの如何に関わらず、最高のスコアを有する終端状態を、復号化したパス中の最終状態として選択する。次いで、各状態で最高の $H_c(g)(t, n)$ スコアを生成するのに使用したエッジおよび t に沿ってトレースバックすることにより、復号化パスを判定する。

【0179】

本発明の下で軌跡のクラスを使用することにより、このようにして復号化を実行することが可能となる。このようなクラスがない場合、状態数が無限となるので、有限状態変換器を使用することができない。その理由は、軌跡 (g) が連続的であり、したがって任意の値を取れるからである。軌跡の代わりに軌跡のクラスについてのスコアを生成することにより、有限の状態数を有することが可能となり、図 3 で述べた復号化が可能となる。各状態に対して計算されるスコアの数制限するためにクラスを使用するが、最適な軌跡 g' は、 g' についての連続的な値の集合から選択され、かつ g' は、スコアを生成するのに使用される（上記の式 41 を参照）ことにも留意されたい。これにより、上記の有限状態復号化が依然として可能となると共に、スコアの判定が正確となる。

【0180】

図 7 に、本発明を使用することができる音声認識システムのブロック図を与える。図 7 では、トレーナまたはユーザのいずれかである話者 700 が、マイクロフォン 704 に向かって発話する。マイクロフォン 704 はまた、1 つまたは複数のノイズ源 702 から加算ノイズも受け取る。マイクロフォン 704 によって検出された可聴信号は電気信号に変換され、アナログ - デジタル (A - D) 変換モジュール 706 に供給される。

【0181】

A - D 変換モジュール 706 は、マイクロフォン 704 からのアナログ信号を一連のデジタル値に変換する。いくつかの実施形態では、A - D 変換モジュール 706 は、16 kHz およびサンプル当り 16 ビットでアナログ信号をサンプリングし、それによって毎秒 32 キロバイトの音声データを生成する。これらのデジタル値はフレームコンストラクタ 707 に供給され、一実施形態ではフレームコンストラクタ 707 は、10 ミリ秒ずれて開始する 25 ミリ秒フレームに値をグループ化する。

10

【0182】

フレームコンストラクタ 707 によって生成されるデータのフレームは特徴抽出モジュール 708 に供給され、特徴抽出モジュール 708 は各フレームから特徴を抽出する。特徴抽出モジュールの例には、線形予測符号化 (LPC) を実施するモジュール、LPC 導出ケプストラム、知覚線形予測 (PLP)、聴覚モデル特徴抽出、およびメル周波数ケプストラム係数 (MFCC) 特徴抽出が含まれる。本発明はこれらの特徴抽出モジュールに限定されず、本発明の状況内で他のモジュールを使用することに留意されたい。

20

【0183】

入力信号がトレーニング信号である場合、この一連の特徴ベクトルはトレーナ 724 に供給され、トレーナ 724 は特徴ベクトルおよびトレーニングテキスト 726 を使用して、本発明の生成モデル 728 をトレーニングする。例えば、前述の EM トレーニングアルゴリズムを使用して、生成モデルをトレーニングすることができる。

【0184】

入力信号がテスト信号である場合、特徴ベクトルがデコーダ 712 に供給され、デコーダ 712 は、特徴ベクトルのストリーム、レキシコン 714、言語モデル 716、および生成モデル 728 に基づいて、最も可能性の高い一続きの単語を識別する。一実施形態の下では、レキシコン (lexicon) 714 は、一続きの特徴ベクトルから単語を識別するために、デコーダ 712 によって走査される有限状態ネットワークを定義する。

30

【0185】

最も可能性の高い一続きの仮説単語が、信頼度測定モジュール 720 に供給される。信頼度測定モジュール 720 は、2 次音響モデル (図示せず) に部分的に基づいて、音声レコグナイザによってどの単語が不適切に識別された可能性が最も高いかを識別する。次いで信頼度測定モジュール 720 は、どの単語が不適切に識別された可能性があるかを示す識別子と共に、一続きの仮説単語を出力モジュール 722 に供給する。本発明を実施するのに信頼度測定モジュール 720 は不要であることを当業者は理解されよう。

【0186】

特定の実施形態を参照しながら本発明を説明したが、本発明の主旨および範囲から逸脱することなく、形態および細部に変更を加えることができることを当業者は理解されよう。

40

【図面の簡単な説明】

【0187】

【図 1】本発明の実施形態の一コンピューティング環境のブロック図である。

【図 2】本発明の実施形態の代替コンピューティング環境のブロック図である。

【図 3】本発明の実施形態の復号化の方法のフローチャートの図である。

【図 4】本発明の実施形態の単純な有限状態図である。

【図 5】本発明の実施形態の様々な状態および時間でのクラスパススコアを示す状態図である。

【図 6】本発明の実施形態の新しい時間インデックスでの状態についてクラスパススコア

50

を計算するのに使用するバス接続を示す状態図である。

【図 7】本発明の実施形態の本発明の一実施形態の下での音声認識システムのブロック図である。

【符号の説明】

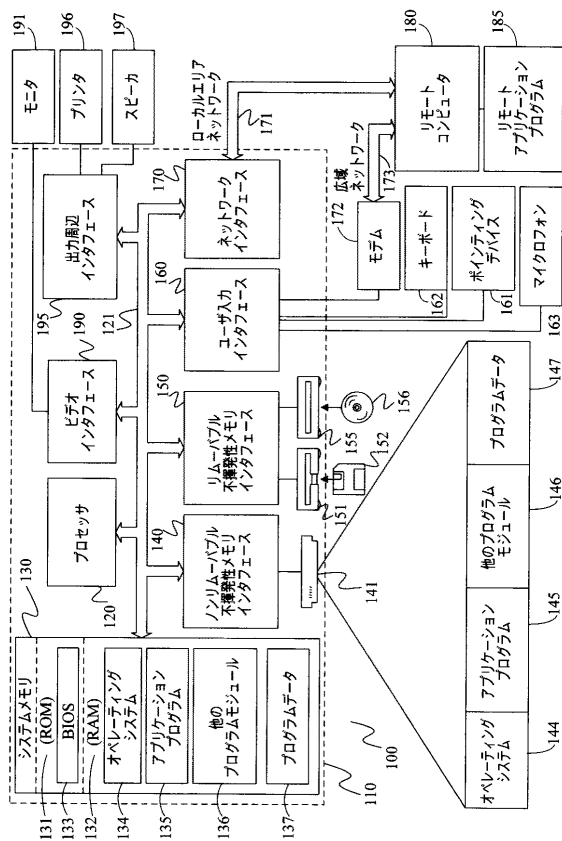
【 0 1 8 8 】

1 0 0	コンピューティングシステム環境	
1 1 0	コンピュータ	
1 2 0	プロセッサ	
1 2 1	システムバス	
1 3 0	システムメモリ	10
1 3 1	読取り専用メモリ (R O M)	
1 3 2	ランダムアクセスメモリ (R A M)	
1 3 3	B I O S	
1 3 4	オペレーティングシステム	
1 3 5	アプリケーションプログラム	
1 3 6	他のプログラムモジュール	
1 3 7	プログラムデータ	
1 4 0	インタフェース	
1 4 1	ハードディスクドライブ	
1 4 4	オペレーティングシステム	20
1 4 5	アプリケーションプログラム	
1 4 6	他のプログラムモジュール	
1 4 7	プログラムデータ	
1 5 0	インタフェース	
1 5 1	磁気ディスクドライブ	
1 5 2	リムーバブル不揮発性磁気ディスク	
1 5 5	光ディスクドライブ	
1 5 6	リムーバブル不揮発性光ディスク	
1 6 0	ユーザ入力インタフェース	
1 7 0	ネットワークインタフェース / アダプタ	30
1 7 1	ローカルエリアネットワーク (L A N)	
1 7 2	モデム	
1 7 3	広域ネットワーク (W A N)	
1 8 0	リモートコンピュータ	
1 8 5	リモートアプリケーションプログラム	
1 9 0	ビデオインタフェース	
1 9 1	モニタ	
1 9 5	出力周辺インタフェース	
1 9 6	プリンタ	
1 9 7	スピーカ	40
2 0 0	モバイル装置	
2 0 2	マイクロプロセッサ	
2 0 4	メモリ	
2 0 6	入出力 (I / O) 構成要素	
2 0 8	通信インタフェース	
2 1 0	バス	
2 1 2	オペレーティングシステム	
2 1 4	アプリケーションプログラム	
2 1 6	オブジェクトストア	
4 0 0、4 0 2、4 0 4、4 0 6	状態	50

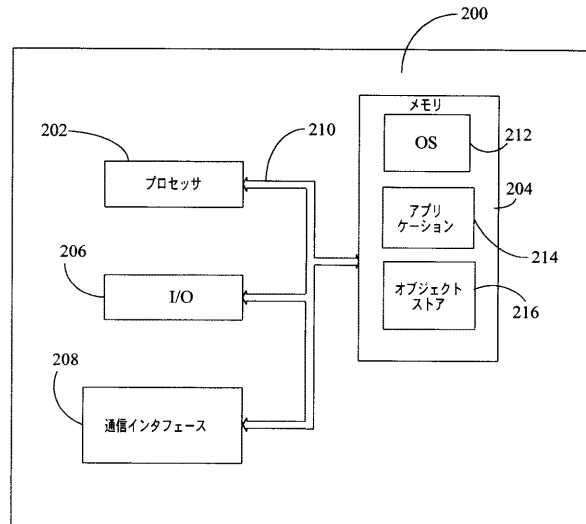
- 700 話者
- 702 ノイズ源
- 704 マイクロフォン
- 706 アナログ - デジタル変換モジュール
- 707 フレームコンストラクタ
- 708 特徴抽出モジュール
- 712 デコーダ
- 714 レキシコン
- 716 言語モデル
- 720 信頼度測定モジュール
- 722 出力モジュール
- 724 トレーナ
- 726 トレーニングテキスト
- 728 生成モデル

10

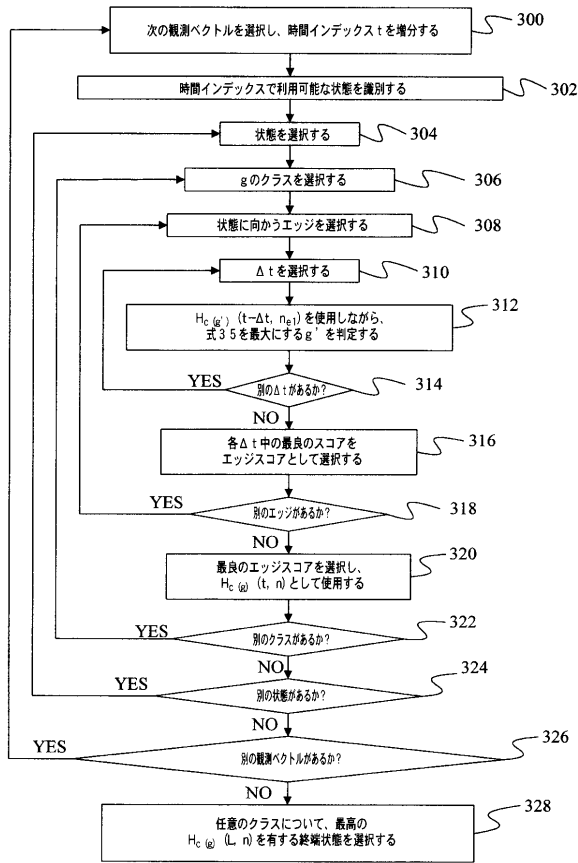
【図1】



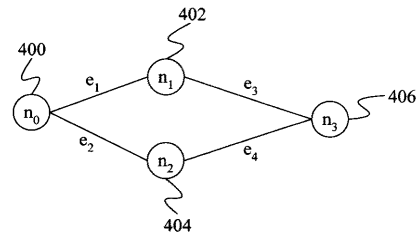
【図2】



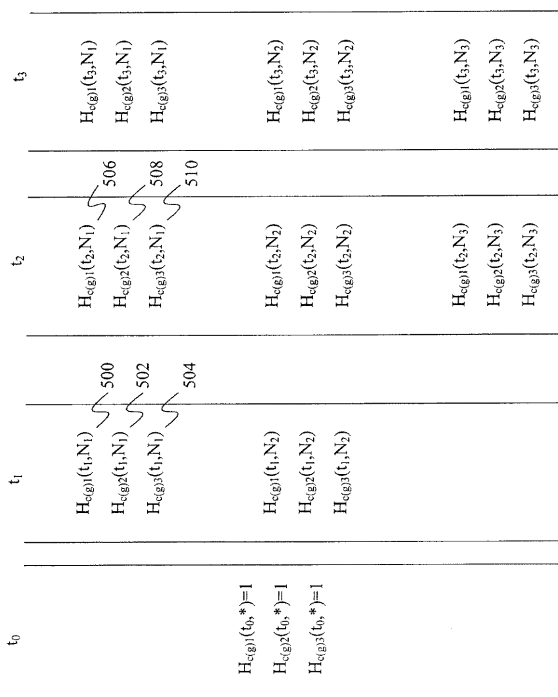
【図 3】



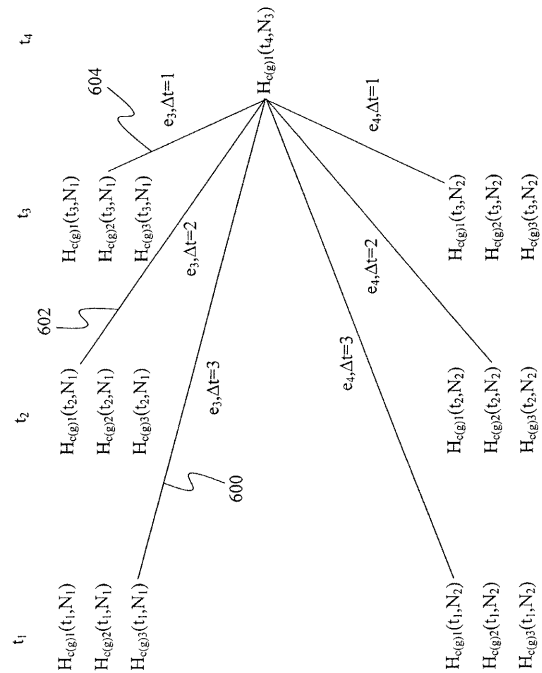
【図 4】



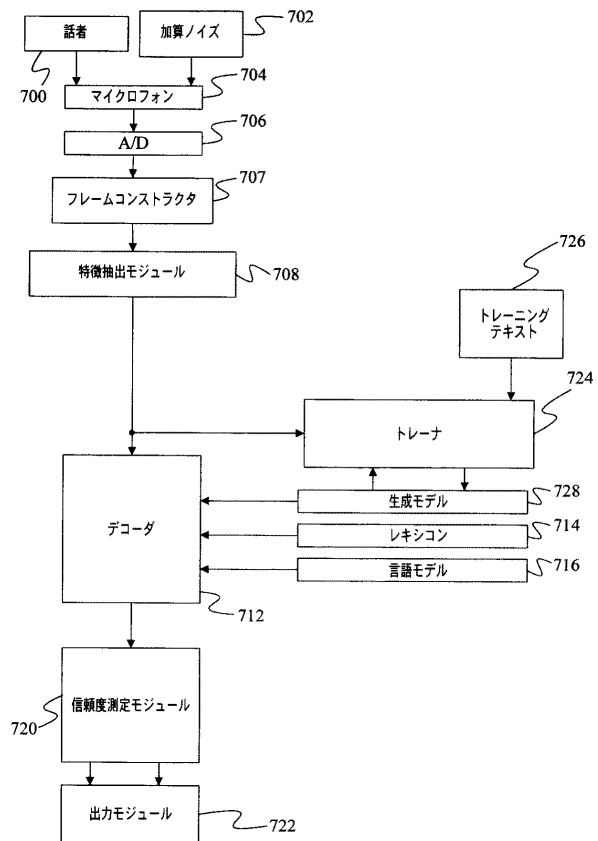
【図 5】



【図 6】



【図 7】



フロントページの続き

- (72)発明者 チョウ ジャン - ライ
中華人民共和国 ペキン チャオヤン ディストリクト シャオヤオジュ ベイリ 113 ビル
ディング 1503
- (72)発明者 フランク トルステン ベルト ザイデ
中華人民共和国 ペキン チャオヤン ディストリクト シャオユン ロード ナンバー32 コ
ンコルディア プラザ アpartment ビー - 1302
- (72)発明者 アセラ ジェイ . アール . グナワルデナ
アメリカ合衆国 98101 ワシントン州 シアトル アラスカン ウェイ 1950 ナンバ
ー323
- (72)発明者 ハガイ アティアス
アメリカ合衆国 98121 ワシントン州 シアトル ウォール ストリート 500 アパー
トメント 901
- (72)発明者 アレハンドロ アセロ
アメリカ合衆国 98006 ワシントン州 ベルビュー 163 プレイス サウスイースト
6525
- (72)発明者 ホアン シュエドン
アメリカ合衆国 98006 - 5345 ワシントン州 ベルビュー 155 アベニュー サウ
スイースト 6179

審査官 山下 剛史

- (56)参考文献 特開平07 - 199379 (JP, A)
特開平08 - 095592 (JP, A)
特開平10 - 063291 (JP, A)
特開平08 - 022296 (JP, A)

- (58)調査した分野(Int.Cl., DB名)
G10L 15/00 - 15/28