

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4528535号
(P4528535)

(45) 発行日 平成22年8月18日(2010.8.18)

(24) 登録日 平成22年6月11日(2010.6.11)

(51) Int. Cl. F I
G 1 0 L 15/06 (2006.01) G 1 0 L 15/06 3 0 0 C
G 1 0 L 15/18 (2006.01) G 1 0 L 15/18 3 0 0 G

請求項の数 10 外国語出願 (全 17 頁)

<p>(21) 出願番号 特願2004-37147 (P2004-37147) (22) 出願日 平成16年2月13日 (2004.2.13) (65) 公開番号 特開2004-246368 (P2004-246368A) (43) 公開日 平成16年9月2日 (2004.9.2) 審査請求日 平成19年1月12日 (2007.1.12) (31) 優先権主張番号 10/365,850 (32) 優先日 平成15年2月13日 (2003.2.13) (33) 優先権主張国 米国 (US)</p>	<p>(73) 特許権者 500046438 マイクロソフト コーポレーション アメリカ合衆国 ワシントン州 9805 2-6399 レッドモンド ワン マイ クロソフト ウェイ (74) 代理人 100077481 弁理士 谷 義一 (74) 代理人 100088915 弁理士 阿部 和夫 (72) 発明者 ミリンド マハジャン アメリカ合衆国 98052 ワシントン 州 レッドモンド ノースイースト 97 ウェイ 17430</p>
---	--

最終頁に続く

(54) 【発明の名称】 テキストから単語誤り率を予測するための方法および装置

(57) 【特許請求の範囲】

【請求項1】

実際の音声単位シーケンスを含む学習テキストから音声信号を生成するステップと、
 生成された音声信号をデコードするステップと、
 デコードされた音声信号を使用して、予測される音声単位シーケンスを生成するステップと、

前記学習テキストに含まれる実際の音声単位シーケンスおよび前記生成された予測される音声単位シーケンスを使用して、音声認識システムのパフォーマンスを識別するためのコンフィュージョンモデルを構築するステップと、

前記コンフィュージョンモデルを使用して、あるテキストが与えられた場合の音声認識システムの誤り率を識別するステップと

を含むことを特徴とする音声認識システムをモデル化する方法。

【請求項2】

請求項1に記載の方法であって、

前記コンフィュージョンモデルおよび言語モデルを使用してテストテキストをデコードするステップと、

前記コンフィュージョンモデルを使用して、前記予測される音声単位シーケンスを表す少なくとも1組のパスを表すテストテキスト用のネットワークを構築するステップと、

構築されたネットワークをスキャンしてパスを探索することで、予測される単語シーケンスを生成するステップと、をさらに備え、

10

20

前記誤り率を識別するステップは、以下の式 1 に従って予測される単語シーケンスに確率を割り当てることを特徴とする方法。

【数 1】

$$p(W_p | W_c) = \frac{P(W_p)}{\sum_{W_p: \varphi(W_p) = \varphi(W_c)} P(W_p)} \cdot \sum_{\bar{t}: \varphi(\bar{t}) = \varphi(W_p)} p(\bar{t} | \varphi_c) \quad \text{式1}$$

ただし、式中 $p(W_p | W_c)$ は、実際の単語シーケンス W_c が与えられた場合の予測される単語シーケンス W_p の確率、 $P(W_p)$ は予測される単語シーケンスの言語モデル確率、分母の総和は、予測される単語シーケンス W_p と同じ音声単位シーケンスを有する予測されるすべての単語シーケンスのすべての言語モデル確率の合計、

10

【数 2】

$$p(\bar{t} | \varphi_c)$$

はネットワークを通るパス

【数 3】

$$\bar{t}$$

に沿った全確率であり、

【数 4】

$$\varphi(W_p)$$

は、単語シーケンス W_p に対応する音声単位シーケンスを表し、

20

【数 5】

$$\varphi(\bar{t})$$

は、パス

【数 6】

$$\bar{t}$$

に沿った予測される音声単位シーケンスを表す。

30

【請求項 3】

請求項 2 に記載の方法であって、

前記誤り率を識別するステップは、予測される単語シーケンス内の単語と、実際の単語シーケンス内の単語との間の差を識別することで、予測される単語シーケンスと実際の単語シーケンスとの間の誤差を決定することを特徴とする方法。

【請求項 4】

請求項 3 に記載の方法であって、

前記誤り率を識別するステップは、前記決定した誤差を以下の式 2 に適用して単語誤り率の期待値を算出することを特徴とする方法。

【数 7】

$$E[\text{WER}] = \frac{\sum_{i=1}^I E[e_i]}{\sum_{i=1}^I N_i} \quad \text{式2}$$

40

ただし、式中 $E[\text{WER}]$ は、テストテキストの単語誤り率の期待値、 $E[e_i]$ はテストテキスト内のシーケンス i での誤差数の期待値、 N_i は実際のシーケンス i 内の単語数、および I はテストテキスト内の文の総数である。

【請求項 5】

前記テストテキストをデコードするステップは、

50

前記第 1 の言語モデルおよび前記コンフュージョンモデルを使用して前記テストテキストをデコードし、予測される第 1 の組の音声単位を生成するステップと、

前記第 1 の言語モデルと異なる第 2 の言語モデルおよび前記コンフュージョンモデルを使用して前記テストテキストをデコードし、予測される第 2 の組の音声単位を生成するステップと、を備え、

前記誤り率を識別するステップは、生成された予測される第 1 および第 2 の組の音声単位を使用して前記誤り率をそれぞれ識別し、識別された誤り率を比較することで、前記第 1 の言語モデルのパフォーマンスを前記第 2 の言語モデルのパフォーマンスと比較するステップを備える

ことを特徴とする請求項 2 に記載の方法。

10

【請求項 6】

前記音声信号をデコードするステップは、学習言語モデルを使用して前記音声信号をデコードすることを含み、前記学習言語モデルは、前記テストテキストのデコードに使用する前記言語モデルと異なる機能を備えることを特徴とする請求項 2 に記載の方法。

【請求項 7】

前記誤り率を使用して前記言語モデルの学習を行うステップをさらに備えることを特徴とする請求項 1 に記載の方法。

【請求項 8】

前記コンフュージョンモデルは、実際の音声単位から少なくとも 1 つの予測される音声単位への変換、およびこうした変換を行う確率をそれぞれ提供する 1 組のルールを含むことを特徴とする請求項 1 に記載の方法。

20

【請求項 9】

前記ルールは、実際のシーケンス内の 1 の音声単位の右および / または左の文脈に依存するルールを含むことを特徴とする請求項 8 に記載の方法。

【請求項 10】

前記コンフュージョンモデルを構築するステップは、学習テキストから生成された音声信号をデコードし、予測される音声単位シーケンスを生成すること、

前記学習テキスト内の学習音声単位のシーケンスを識別すること、

前記予測される音声単位シーケンスを前記学習音声単位のシーケンスに合わせて配列すること、および

30

前記配列されたシーケンスを使用して前記コンフュージョンモデルの学習を行うことを含むことを特徴とする請求項 1 に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は音声認識に関する。より詳細には、本発明は音声認識システムのパフォーマンスのモデル化に関する。

【背景技術】

【0002】

40

音声認識では、音響モデルおよび言語モデルを使用して音響信号を単語シーケンスに変換する。音響モデルは、音響信号の特徴を、確率を備える単音などの可能なサブワード音声単位 (sub-word speech unit) シーケンスに変換する。言語モデルは、音響モデルによって識別された音のシーケンスから形成することができる様々な単語シーケンスの確率分布を提供する。

【0003】

音響モデルの学習は、一般に、話者に既知のテキストを読ませ、次いで学習音声から学習テキストを予測するように音響モデルを作り上げることによって行われる。言語モデルの学習は、一般に、テキストの大きいコーパスから、コーパス内の様々な単語シーケンスの確率を単に識別することによって行われる。

50

【発明の開示】

【発明が解決しようとする課題】

【0004】

得られる音声認識システムのパフォーマンスは、音響モデルおよび言語モデルの学習に使用される学習テキストにある程度依存する。その結果、音声認識システムは、一部のタスクドメインにおいて、他のタスクドメインにおけるより良好に機能する。音声認識システムが特定のタスクドメインでどのように動作するかを決定するために、タスクの実行時にユーザが使用されると思われる単語を誰かが話して、システムによってデコード可能な音響データを生成する必要がある。音声認識システムのパフォーマンスを決定するのに十分な量の音響データを生成するために人を雇うことは、費用がかかり、音声で実行可能になるコンピュータアプリケーションの開発の障壁となっている。

10

【0005】

さらに、音響データを生成することは費用がかかるため、こうしたデータは、言語モデルの学習に使用するコーパス全体については生成されていなかった。その結果、言語モデルの学習は一般に、音響モデルが言語モデルコーパスに対してどのように機能するかを検査することなく行われていた。したがって、音響データを必要とすることなく音響モデルと言語モデルの組合せのパフォーマンスの測定にテキストのコーパスを使用できるようにするシステムを得ることが有益となる。これによって、音響モデルとの組合せで言語モデルの識別学習 (discriminative training) が可能となる。

20

【課題を解決するための手段】

【0006】

音声認識システムをモデル化する方法は、学習テキストから生成された音声信号をデコードして予測される音声単位シーケンスを生成することを含む。学習テキストは実際の音声単位シーケンスを含んでおり、これを予測される音声単位シーケンスとともに使用してコンフュージョンモデル (confusion model) が形成される。別の実施形態では、コンフュージョンモデルを使用してテキストをデコードして、音声認識システムがテキストに基づいて音声信号をデコードした場合に予想される誤り率を識別する。

【発明を実施するための最良の形態】

【0007】

図1は、本発明を実施するのに適したコンピューティングシステム環境100の例を示している。コンピューティングシステム環境100は、適したコンピューティング環境の一例にすぎず、本発明の使用または機能の範囲に関する限定を示唆するものではない。また、コンピューティング環境100を、動作環境100の例に示した構成要素のいずれか1つ、またはその組合せに関連する依存性または必要条件を有しているものと解釈すべきではない。

30

【0008】

本発明は、他の多くの汎用または専用コンピューティングシステム環境または構成で動作可能である。本発明との使用に適したよく知られている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、それだけには限定されないが、パーソナルコンピュータ、サーバーコンピュータ、ハンドヘルドまたはラップトップ装置、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、プログラム可能な家庭用電化製品、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、テレフォニーシステム、上記の任意のシステムまたは装置を含む分散コンピューティング環境などがある。

40

【0009】

本発明は、コンピュータによって実行される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な文脈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する、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、構成要素、データ構造などを含む。本発明は、タスクが通信ネットワークによってリンクされているリモート処理装置によって実行される分散コンピューティング

50

環境で実施するように設計されている。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を、メモリ記憶装置を含むローカルおよびリモートのコンピュータ記憶媒体に置くことができる。

【0010】

図1を参照すると、本発明を実施するシステムの例は、汎用コンピューティング装置をコンピュータ110の形で含んでいる。コンピュータ110の構成要素は、それだけには限定されないが、処理装置120、システムメモリ130、およびシステムメモリを含む様々なシステム構成要素を処理装置120に結合するシステムバス121を含む。システムバス121は、様々なバスアーキテクチャのうちの任意のものを使用するメモリバスまたはメモリコントローラ、周辺機器バス、およびローカルバスを含むいくつかのタイプのバス構造のうちどんなものでもよい。こうしたアーキテクチャには、それだけには限定されないが一例として、業界標準アーキテクチャ（ISA）バス、マイクロチャンネルアーキテクチャ（MCA）バス、拡張ISA（EISA）バス、ビデオ電子装置規格化協会（VESA）ローカルバス、およびメザンバスとしても知られている周辺機器相互接続（PCI）バスなどがある。

10

【0011】

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からアクセスできる使用可能な任意の媒体とすることができ、揮発性および不揮発性媒体、リムーバブルおよび非リムーバブル媒体を含む。コンピュータ可読媒体は、それだけには限定されないが一例として、コンピュータ記憶媒体および通信媒体を含み得る。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、他のデータなど、情報を記憶するための任意の方法または技術で実施される揮発性および不揮発性のリムーバブルおよび非リムーバブル媒体がある。コンピュータ記憶媒体には、それだけには限定されないが、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、デジタル多用途ディスク（DVD）または他の光ディスク記憶装置、磁気カセット、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、または所望の情報の格納に使用でき、コンピュータ110からアクセスできる他の任意の媒体などがある。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変調されたデータ信号に組み込む。これには任意の情報配送媒体がある。「変調されたデータ信号」という用語は、1つまたは複数のその特徴が、信号内の情報を符号化するように設定または変更された信号を意味する。通信媒体には、それだけには限定されないが一例として、有線ネットワーク、直接配線された接続などの有線媒体、および音響、RF、赤外線、その他の無線媒体などの無線媒体がある。また、上記のどんな組合せでもコンピュータ可読媒体の範囲内に含まれるものとする。

20

30

【0012】

システムメモリ130は、読取り専用メモリ（ROM）131やランダムアクセスメモリ（RAM）132など、揮発性および/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。基本入出力システム133（BIOS）は、例えば起動中など、コンピュータ110内の要素間での情報の転送を助ける基本ルーチンを含み、一般にROM131に格納されている。RAM132は一般に、処理装置120から直接アクセス可能な、かつ/または処理装置120が現在処理しているデータおよび/またはプログラムモジュールを含む。図1は、それだけには限定されないが一例として、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を示している。

40

【0013】

コンピュータ110は、他のリムーバブル/非リムーバブル、揮発性/不揮発性コンピュータ記憶媒体を含むこともできる。一例にすぎないが、図1は、非リムーバブル不揮発性磁気媒体から読み取り、あるいはそこに書き込むハードディスクドライブ141、リムーバブル不揮発性磁気ディスク152から読み取り、あるいはそこに書き込む磁気ディス

50

クドライブ151、およびCD-ROMや他の光媒体など、リムーバブル不揮発性光ディスク156から読み取り、あるいはそこに書き込む光ディスクドライブ155を示している。動作環境の例で使用できる他のリムーバブル/非リムーバブル、揮発性/不揮発性コンピュータ記憶媒体には、それだけには限定されないが、磁気テープカセット、フラッシュメモリカード、デジタル多用途ディスク、デジタルビデオテープ、半導体RAM、半導体ROMなどがある。ハードディスクドライブ141は一般に、インターフェイス140などの非リムーバブルメモリインターフェイ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、インターフェイス150などのリムーバブルメモリインターフェイスによってシステムバス121に接続される。

10

【0014】

上述し、図1に示したドライブおよびその関連のコンピュータ記憶媒体は、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、およびコンピュータ110の他のデータの記憶を提供する。図1では例えば、ハードディスクドライブ141は、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を記憶するものとして示されている。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注意されたい。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は少

20

【0015】

ユーザは、キーボード162、マイクロフォン163、および一般にマウス、トラックボール、タッチパッドなどのポインティング装置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、衛星放送受信アンテナ、スキャナなどがある。これらおよび他の入力装置は、しばしばシステムバスに結合されているユーザ入力インターフェイス160を介して処理装置120に接続されるが、パラレルポート、ゲームポート、ユニバーサルシリアルバス(USB)など他のインターフェイスおよびバス構造で接続

30

【0016】

コンピュータ110は、リモートコンピュータ180など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操作される。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、サーバー、ルーター、ネットワークPC、ピア装置、または他の一般のネットワークノードでよく、一般にコンピュータ110に関連して上述した多くまたはすべての要素を含む。図1に示した論理接続は、ローカルエリアネットワーク(LAN)171および広域エリアネットワーク(WAN)173を含むが、他のネットワークを含んでいてもよい。こうしたネットワーキング環境は、オフィス、全社規模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はごく一般的である。

40

【0017】

LANネットワーキング環境で使用する場合、コンピュータ110は、ネットワークインターフェイスまたはアダプタ170を介してLAN171に接続される。WANネットワーキング環境で使用する場合、コンピュータ110は一般に、モデム172、またはインターネットなどWAN173にわたって通信を確立する他の手段を含む。モデム172は、内蔵のものでも外付けのものでもよく、ユーザ入力インターフェイス160または他

50

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。ネットワーク式環境では、コンピュータ110に関連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に格納することができる。図1は、それだけには限定されないが一例として、リモートアプリケーションプログラム185をメモリコンピュータ180上に存在するものとして示している。図示したネットワーク接続は例であり、コンピュータ間の通信リンクを確立する他の手段を使用してもよいことは理解されよう。

【0018】

本発明では、音声認識システムのパフォーマンスのモデルを構築し、それを使用して、あるテキストが与えられた場合の音声認識システムの見込み誤り率(likely error rate)を識別する。以下の説明で、このモデルをコンフュージョンモデルと呼ぶ。コンフュージョンモデルを使用することによって、特定のテキストから生成された音声をデコードする際に音声認識システムがどのように機能するかを推定するのに音響データは必要なくなる。

10

【0019】

図2は、本発明によるコンフュージョンモデルを構築し、使用方法を示すフロー図である。図2のステップ200で、音響モデル300の学習が行われる。これには、話者308に学習テキスト304の一部を朗読させて音響信号を生成することが必要となる。この音響信号は受信機309によって検出され、電子信号に変換される。電子信号は、特徴抽出機に提供され、そこで電子信号から1つまたは複数の特徴が抽出される。こうした特徴ベクトルを識別する方法は、当分野ではよく知られており、これには、特徴ベクトルごとに13のケプストラム値を生成する13次元メル周波数ケプストラム係数(MFCC: Mel Frequency Cepstrum Coefficients)抽出などがある。ケプストラム特徴ベクトルは、電子信号の対応するフレーム内にある音声信号のスペクトルの内容を表す。抽出された特徴ベクトルは、トレーナ302に提供され、トレーナはその特徴を使用して音響モデル300の学習を行う。音響モデルを構築する方法は、当分野ではよく知られている。音響モデルは任意の所望の形態でよく、これには、それだけには限定されないが、隠れマルコフモデルなどがある。

20

【0020】

ステップ202で、話者308が学習テキスト304の一部を話して、学習が行われた音響モデルを使用してデコードされるテスト信号が生成される。音響信号は、受信機309および特徴抽出機310によって特徴ベクトルに変換され、特徴ベクトルは、デコーダ312に提供される。

30

【0021】

デコーダ312は、音響モデル300、辞書315、および学習言語モデル314を使用して、特徴を、予測される音声単位シーケンスにデコードする。多くの実施形態では、音声単位とは、単音、2音、3音など音の単位(phonetic unit)である。

【0022】

ステップ200および202は、1個抜き(leave-one-out)手法を使用して行うことができ、学習テキスト304が複数の組の学習データに分割されることに注意されたい。1組を除く全部を使用して音響モデルが構築され、モデルの構築に使用されない組は、音響モデルおよび言語モデルを使用してデコードされる。次いでステップ200および202が繰り返され、デコードすべき組として、学習データから別の1組が選択され、残りの組から音響モデルが構築される。これは学習テキスト304内のデータの組ごとに繰り返され、その結果学習テキスト304内のデータの組ごとに個々の組の予測される音声単位が提供される。

40

【0023】

学習データのデコードに使用する言語モデルの学習が学習データを使用して行われていれば、同様の1個抜き手順を使用して言語モデルの学習を行う必要がある。1個抜き手法は、予測される音声単位における偏りを避けるのに有益である。

【0024】

50

ステップ204で、予測される音声単位シーケンスが、学習テキスト304からの実際の音声単位シーケンスに合わせて配列される。これは、図3の配列モジュール316によって行われる。一実施形態では、この配列は動的プログラミングを使用して行われ、可能な配列が互いに比較され、ある目的関数の最適化に基づいて可能な配列の最適なものが選択される。一実施形態では、この目的関数は、配列が終了した後で予測される音声単位シーケンスが実際の音声単位シーケンスと異なる度合いを示す誤差関数である。一部の実施形態では、誤差関数は簡単な二値関数であり、並べられた2つの音声単位が一致しない場合は誤差値1が生成され、2つの音声単位が一致する場合は誤差値0が生成される。別の実施形態では、異なるタイプの誤差に異なる重みが付加される。例えば、似た音を提供する音声単位は、大幅に異なる音を提供する音声単位より誤差値が小さい可能性がある。

10

【0025】

予測されるシーケンスを実際の音声単位シーケンスに合わせて配列するステップは、コンフュージョンモデルが、学習が可能となる前に、並べられた音声単位を必要とする場合にのみ、ステップ204で行われる。隠れマルコフモデルなど、いくつかのタイプのコンフュージョンモデルでは、こうした配列は不要である。というのは、モデルの学習には、音声単位を配列するステップが本来含まれているからである。

【0026】

ステップ206で、コンフュージョンモデルトレーナ318によってコンフュージョンモデル320が構築される。一実施形態では、コンフュージョンモデル320として隠れマルコフモデルが構築される。

20

【0027】

図4は、本発明のコンフュージョンモデルの一実施形態による4状態隠れマルコフモデルを示す。本発明のこの実施形態では、個々の隠れマルコフモデルが実際の音声単位ごとに構築される。各隠れマルコフモデルは、その特定の実際の音声単位から生成できる予測される音声単位を出力として生成する。

【0028】

4状態モデルには、状態400、402、404、および406がある。コンフュージョンモデルでは、状態400および406は、非生成状態(non-emitting state)であり、こうした状態を離れる遷移からは、予測される音声単位は生成されない。一方、状態402および404を離れる遷移からは、こうした遷移が使用されるたびに予測される音声単位が1つ生成される。状態402および404では、可能な予測される各音声単位を生成する個々の出力確率が存在する。

30

【0029】

これらの状態は、関連する確率を有する遷移によって連結される。ある状態を離れる遷移の確率の合計は1に等しい。したがって、遷移408および410の遷移確率の合計は1に等しく、また遷移412および414の確率の合計は1に等しい。

【0030】

図4のモデルを通る様々なパスは、様々な置換および挿入の可能性を表す。例えば、遷移410に沿った状態400から状態406へのパスは、実際の音声単位の代わりに置き換えられる音声単位が予測される音声単位シーケンス内にないことを表す。遷移408および414に沿った、状態400から402を通る406までのパスは、実際の音声単位が予測される音声単位1つに置き換えられることを表す。このパスは、実際の音声単位が予測される音声単位と同じである状況を含む。遷移408、412、および418に沿った、状態400から状態402および404を通る406までのパスは、実際の音声単位が予測される音声単位2つに置き換えられることを表す。予測される音声単位の一つは状態402からの遷移で生成され、第2の音声単位は状態404からの遷移で生成される。

40

【0031】

状態404は、実際の音声単位の代わりに任意の数の予測される音声単位を得る機構を提供する自己遷移416を含む。

【0032】

50

図4のモデルなどの各隠れマルコフモデルの学習は、まず、状態402および404における予測される音声単位ごとに等しい確率を想定し、遷移ごとに等しい遷移確率を想定することによって行われる。この簡単なモデルを使用して、実際の音声単位と予測される音声単位との最適な配列を識別し、それによってHMMごとに状態確率および遷移確率の学習に使用できる配列が提供される。更新されたモデルで、配列が調整され、新しい配列を使用してモデルの再学習が行われる。これはモデルが安定するまで続けられる。この学習手順は、ビタビ学習としても知られている。この学習プロセスのわずかに異なるバリエーションでは、実際の音声単位と予測される音声単位との複数の配列を考えることができる。これは、Forward-Backward学習またはBaum-Welch学習として知られている。

10

【0033】

他の実施形態では、コンフュージョンモデルは、各ルールが実際の音声単位シーケンス内の単一の音声単位から予測される音声単位シーケンス内の0個、1個、または複数個の音声単位への変換の可能性を提供するルールベースのモデルとして構築される。本発明の実施形態では、各ルールは、文脈とは無関係とすることも文脈に依存することもできる。文脈依存型ルールでは、ルールは、実際のシーケンス内の単一の音声単位の左、右、または左右の文脈に依存することができる。さらに、本発明では、単一の音声単位の左または右の任意の数の音声単位を使用することができ、異なる文脈長さの組合せをいっしょに使用して、一方のルールが一方の文脈長さを使用し、もう一方のルールがより長い文脈長さを使用する状態で、異なる2つのルールが実際の音声単位を予測される音声単位シーケンスに変換する異なる2つの確率を提供できるようにすることができる。内挿法により確率を組み合わせる、または文脈に基づいて単一のルール確率を選択することができる。

20

【0034】

ルール確率の学習を行うには、実際の音声単位シーケンスと予測される音声単位シーケンスの間の配列が検査されて、実際の音声単位を予測される音声単位シーケンスに変換するために各ルールを使用できる回数を決定する。この数を、実際の音声単位シーケンス内の単一の音声単位がルールに関連する特定の文脈内に見つかる回数で割る。したがって、各確率は、実際の音声単位、および必要に応じて実際の音声単位シーケンス内の特定の文脈が与えられた場合に予測される音声単位シーケンスを生成する尤度を示す。

【0035】

本発明の一実施形態では、実際の音声単位の置換を必要とすることなく音声単位を予測されるシーケンスに挿入することができる。これは、配列前の実際の音声単位シーケンス内の各音声単位の間空の音声単位を挿入することによって達成される。配列中、予測されるシーケンス内には、こうした空の音声単位の相手として配列されるものがない。しかし、音声単位は、時々、予測されるシーケンス内の1つまたは複数の音声単位に合わせて配列されることがある。その結果、空の音声単位から配列内で見つかる予測される音声単位に変換するルールが生成される。

30

【0036】

本発明の一実施形態では、ルールごとの確率が生成された後、一部のルールが削除されてルールセット内のルールの数を低減させ、それによってデコードが簡素化される。この削除は、文脈および単一の音声単位が実際の音声単位シーケンス内で見つかる回数に基づいて、または文脈および実際の音声単位が予測される音声単位の特定のシーケンスを生成する回数に基づいて行うことができる。後者の場合、削除されたが、別の予測される音声単位シーケンスを提供したルールと同じ文脈を共有するルールの確率を、同じ文脈および実際の音声単位での確率の合計が1に等しくなるように再計算する必要がある。

40

【0037】

ステップ206でコンフュージョンモデルが構築された後、これを使用して、音響信号の生成の必要なく音響モデルのパフォーマンスをモデル化することができる。ステップ208で、デコーダ502は、コンフュージョンモデル504、辞書508、および言語モデル506を使用して、図5のテストテキスト500をデコードする。デコーダ502は

50

、1組の予測される単語シーケンスを、コンフュージョンモデル504、辞書508、および言語モデル506を使用して計算されたシーケンスごとの確率とともに生成する。また、デコーダ502は、テストテキスト500内の実際の単語シーケンスを渡す。

【0038】

デコーダ502が使用する言語モデル506は、ほとんどの実施形態では、学習言語モデル314とは異なる。一般に、学習言語モデル314は脆弱な言語モデルとなるように選択されており、したがってコンフュージョンモデルは、音響モデル300の弱点を示す。一方、言語モデル506は、音声認識システムで使用される言語モデルとより類似した強力な言語モデルである。一部の実施形態では、異なる言語モデルを同じコンフュージョンモデル504およびテストテキスト500とともに使用して、算出された単語誤り率に基づいて2つの言語モデルの相対的パフォーマンスを決定する。これについて以下でさらに詳しく論じる。

10

【0039】

ルールベースのコンフュージョンモデルでは、デコーダ502は、ルールをコンフュージョンモデル504に適用して、図6のパスなど、可能な予測される音声単位シーケンスを表す1組のパスを表すネットワークを構築する。図6には、状態600、602、604、606、608、および610の間に遷移を示している。任意の2つの状態の間の遷移は、テストテキスト500で表される実際の音声単位シーケンス内の特定の音声単位の代わりに使用されるある確率を有する予測される音声単位シーケンスを表す。例えば、状態602と604の間の遷移は、音声単位P1の代わりに使用できる予測されるすべての音声単位シーケンスを表す。同様に、状態604と606の間の遷移は、実際の音声単位シーケンスの音声単位P1と音声単位P2の間に挿入することができるすべての遷移を表す。遷移は、予測される音声単位シーケンス内に音声単位がない場合も提供することに注意されたい。したがって、状態602と604の間の遷移の1つは、音声単位がないことを表すを含み得る。これは、予測される音声単位シーケンス内から実際の音声単位P1が単に除去されることを示す。ネットワークは、空の音声単位に対応する状態および遷移も含む。これらは、図6に示すように、状態600と602、状態604と606の間に、実際の音声単位とともに交互の位置に挿入される。これらの遷移は、空の音声単位に対応する挿入ルールに対応する。

20

【0040】

図6では、各遷移は特定の確率を有している。このテストテキスト用に生成されたネットワークから予測される単語シーケンスを生成するために、デコーダは、ネットワークをスキャンすることによってネットワーク内のパスを探索する。ここでは、深さ優先探索や幅優先探索などの標準グラフスキャン戦略(graph traversal strategies)を使用することができる。デコーダが探索した各パスは、予測される音声単位シーケンスに対応する。

30

【0041】

効率を高め、デコーダが探索するパスの数を制限するために、デコーダは、スキャン中に生成された部分パスを削除する(すなわちそれ以上の考察から除外する)ことを選択することができる。例えば、部分パスが、任意の単語シーケンスに対応する任意の音声単位シーケンスの先頭部分に一致していない予測される音声単位シーケンスに対応している場合、こうした部分パスを削除することができる。他の部分パスまたは完全パスに比べて低い確率を有する部分パスも削除することができる。

40

【0042】

削除についてさらに説明するために、デコーダは、まず実際の音声単位シーケンスと同じ予測される音声単位シーケンスを生成してもよい。こうしたシーケンスは一般に、高い確率を有しており、低い確率の他のパスの識別を助ける。

【0043】

デコーダ502は、特定のパスに沿って終了状態612に到達するとき、探索されたパスに対応する予測される音声単位シーケンスを、その関連する確率とともに使用すること

50

ができる。次いでデコーダは、単語単位シーケンスに対応する音声単位シーケンスを予測される音声単位シーケンスと照合することによって、予測される音声単位シーケンスを生成することができるすべての単語シーケンスを識別する。

【 0 0 4 4 】

デコーダは、ネットワーク内のすべてのパスを探索した後、下記の式に従って予測される単語シーケンスに確率を割り当てる。説明を簡潔にするために、式の導出の際に単純化するための仮説をいくつか立てたことに注意されたい。例えば、単語ごとに単一の発音（すなわち1つの音声単位シーケンス）があると仮定した。

【 0 0 4 5 】

【数 1】

$$p(W_p | W_c) = \frac{P(W_p)}{\sum_{W_p: \phi(W_p) = \phi(W_c)} P(W_p)} \cdot \sum_{\bar{t}: \phi(\bar{t}) = \phi(W_p)} p(\bar{t} | \phi_c) \quad \text{式1}$$

10

【 0 0 4 6 】

式中 $p(W_p | W_c)$ は、実際の単語シーケンス W_c が与えられた場合の予測される単語シーケンス W_p の確率、 $P(W_p)$ は予測される単語シーケンスの言語モデル確率、分母の総和は、予測される単語シーケンス W_p と同じ音声単位シーケンスを有する予測されるすべての単語シーケンスのすべての言語モデル確率の合計、および

【 0 0 4 7 】

【数 2】

$$p(\bar{t} | \phi_c)$$

【 0 0 4 8 】

はネットワークを通るパス

【 0 0 4 9 】

【数 3】

$$\bar{t}$$

【 0 0 5 0 】

に沿った全確率であり、 W_p のシーケンスと同じ音声単位シーケンスを生成するすべてのパスについて和を取る。

【 0 0 5 1 】

【数 4】

$$\phi(W_p)$$

【 0 0 5 2 】

は、単語シーケンス W_p に対応する音声単位シーケンスを表し、

【 0 0 5 3 】

【数 5】

$$\phi(\bar{t})$$

【 0 0 5 4 】

は、パス

【 0 0 5 5 】

【数 6】

$$\bar{t}$$

【 0 0 5 6 】

に沿った予測される音声単位シーケンスを表す。パスごとの確率

20

30

40

50

【 0 0 5 7 】

【 数 7 】

 $p(\bar{t}|\varphi_c)$

【 0 0 5 8 】

は、そのパスに沿った各遷移に関連する個々の確率の積として求められる。デコーダのトレリスを通る異なるパスは、遷移のために、同じ予測される音声単位シーケンスを有し得ることに注意されたい。遷移は、パスの類似性を決定するために音声単位として見なされない。例えば、シーケンス $t - \quad - i y$ は、シーケンス $t - i y - \quad$ と同じであると見なされることになる。というのは、シーケンスの類似性を決定する際に音声単位は無視されるからである。

10

【 0 0 5 9 】

予測される単語シーケンスごとの確率が決定された後、ステップ 2 1 0 で、予測される単語シーケンスと実際の単語シーケンスの間の誤差が単語誤り率計算機 5 1 0 によって識別される。誤差は、予測される単語シーケンス内の単語と、デコーダ 5 0 2 によって提供される実際の単語シーケンス内の単語の間の差を識別することによって決定される。

【 0 0 6 0 】

ステップ 2 1 2 で、誤差を使用してテストテキスト 5 0 0 の単語誤り率を生成する。一実施形態では、単語誤り率の期待値は、次のように決定される。

【 0 0 6 1 】

【 数 8 】

$$E[\text{WER}] = \frac{\sum_{i=1}^I E[e_i]}{\sum_{i=1}^I N_i}$$

式 2

20

【 0 0 6 2 】

式中 $E[\text{WER}]$ は、テストテキスト 5 0 0 の単語誤り率の期待値、 $E[e_i]$ はテストテキスト内のシーケンス i での誤差数の期待値、 N_i は実際のシーケンス i 内の単語数、および I はテストテキスト 5 0 0 内の文の総数である。

30

【 0 0 6 3 】

一実施形態では、 $E[e_i]$ は文 i について予測される単語シーケンスごとの誤差数と、その予測される単語シーケンスの確率の積の合計に等しい。これは、次のような式で書き表すこともできる。

【 0 0 6 4 】

【 数 9 】

$$E[e_i] = \sum_{w_p} \text{ErrCount}(w_p, W_c) * p(w_p | W_c)$$

式 3

【 0 0 6 5 】

ここで、 W_c はシーケンス i 内の実際の単語シーケンス、 $\text{ErrCount}(w_p, W_c)$ は予測される単語シーケンス w_p を実際の単語シーケンス W_c と照合することによって識別される誤差数、および $p(w_p | W_c)$ は実際の単語シーケンス W_c が与えられた場合の予測される単語シーケンス w_p の確率である。予測されるすべての単語シーケンスについて和を取る。

40

【 0 0 6 6 】

他の実施形態では、ある文の予想される誤差数は、生じた誤差のタイプに基づいて誤差に重み付けすることによって計算される。言い換えれば、 $\text{ErrCount}(w_p, W_c)$ の計算は、誤差のタイプに基づいて誤差に重み付けすることによって行われる。

【 0 0 6 7 】

50

したがって、式2を使用してテストテキストの単語誤り率を生成することが可能となる。テストテキスト内の各文は別々に検査されるため、音声認識システムがデコードするときに多数の誤差を生成しそうなテストテキスト内の単語シーケンスを識別することも可能である。上に示したシステムでは、音声認識システムがそのテキストに対してどのように機能するかを決定するのに音響データは必要ないことに注意されたい。これによって、様々なタスクに関係する音声認識システムの評価コストが大幅に低減される。また、これによって開発者は、ユーザに入力を要求し、その結果音声認識システムがより容易にデコードできる単語シーケンスを使用するようユーザを導くやり方を変えることができる。

【0068】

さらに、本発明によって複数の言語モデルを互いに比較することができるようになる。これは、1つの言語モデルを使用して単語誤り率を決定し、次いで第2の言語モデルを使用して誤り率を決定することによって行うことができる。次いで単語誤り率が互いに比較されて言語モデルの相対的なパフォーマンスが決定される。

10

【0069】

さらに本発明によって、言語モデルの学習を、ステップ212で計算された単語誤り率に基づく識別学習を使用して音響モデルで行うことができる。こうした学習で、言語モデルは、単語誤り率を改善するように変更される。この学習は一部音響モデルのパフォーマンスに基づいているため、結果として得られる言語モデルは、音響モデルのパフォーマンスを参照することなしに学習が行われる言語モデルより良好に機能すると考えられる。

【0070】

20

本発明は、特定の実施形態に関連して説明してきたが、本発明の意図および範囲から逸脱することなく、形態および詳細に変更を加えることができることを当分野の技術者であれば理解されよう。

【図面の簡単な説明】

【0071】

【図1】本発明の実施形態を実施できる一般のコンピューティング環境を示すブロック図である。

【図2】本発明の実施形態によるコンフュージョンモデルの構築および使用の方法を示すフロー図である。

【図3】本発明の実施形態によるコンフュージョンモデルの学習に使用する構成要素を示すブロック図である。

30

【図4】本発明の一実施形態によるHMMコンフュージョンモデルを示す状態図である。

【図5】テキストおよびコンフュージョンモデルを使用して単語誤り率を決定するデコーダを示すブロック図である。

【図6】本発明の一実施形態によるデコード中に形成されたパスを示すトレリス図である。

【符号の説明】

【0072】

- 100 コンピューティングシステム環境
- 110 コンピュータ
- 120 処理装置
- 121 システムバス
- 130 システムメモリ
- 131 読取り専用メモリ (ROM)
- 132 ランダムアクセスメモリ (RAM)
- 133 基本入出力システム (BIOS)
- 134 オペレーティングシステム
- 135 アプリケーションプログラム
- 136 他のプログラムモジュール
- 137 プログラムデータ

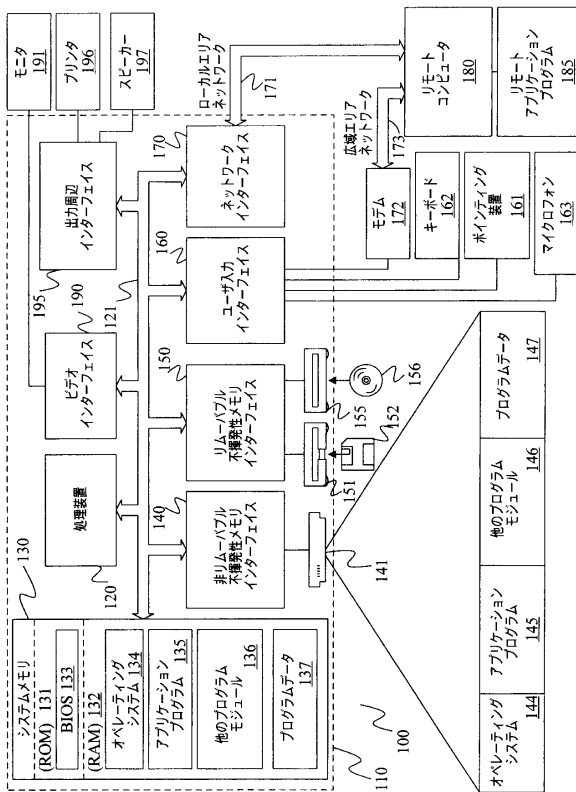
40

50

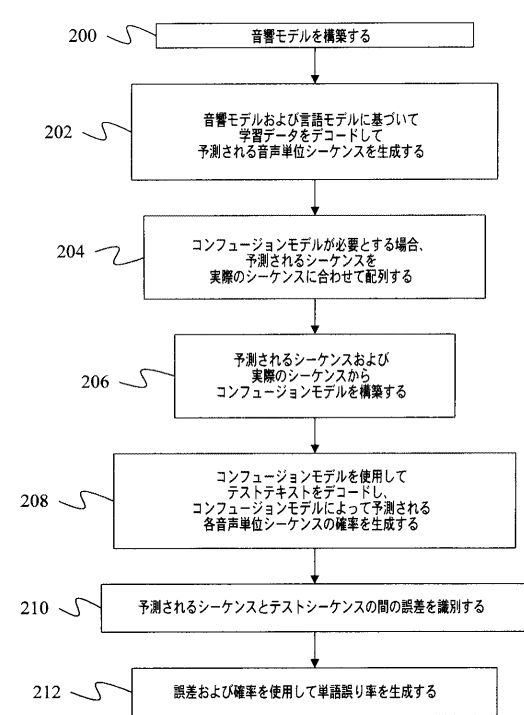
1 4 0	インターフェイス	
1 4 1	ハードディスクドライブ	
1 4 4	オペレーティングシステム	
1 4 5	アプリケーションプログラム	
1 4 6	他のプログラムモジュール	
1 4 7	プログラムデータ	
1 5 0	インターフェイス	
1 5 1	磁気ディスクドライブ	
1 5 2	リムーバブル不揮発性磁気ディスク	
1 5 5	光ディスクドライブ	10
1 5 6	リムーバブル不揮発性光ディスク	
1 6 0	ユーザ入力インターフェイス	
1 6 1	ポインティング装置	
1 6 2	キーボード	
1 6 3	マイクロフォン	
1 7 0	ネットワークインターフェイスまたはアダプタ	
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	広域エリアネットワーク (W A N)	
1 8 0	リモートコンピュータ	20
1 8 5	リモートアプリケーションプログラム	
1 9 0	ビデオインターフェイス	
1 9 1	モニタ	
1 9 5	出力周辺インターフェイス	
1 9 6	プリンタ	
1 9 7	スピーカー	
3 0 0	音響モデル	
3 0 2	トレーナ	
3 0 4	学習テキスト	
3 0 8	話者	30
3 0 9	受信機	
3 1 0	特徴抽出機	
3 1 2	デコーダ	
3 1 4	学習言語モデル	
3 1 5	辞書	
3 1 6	配列モジュール	
3 1 8	コンフュージョンモデルトレーナ	
3 2 0	コンフュージョンモデル	
4 0 0	状態	
4 0 2	状態	40
4 0 4	状態	
4 0 6	状態	
4 0 8	遷移	
4 1 0	遷移	
4 1 2	遷移	
4 1 4	遷移	
5 0 0	テストテキスト	
5 0 2	デコーダ	
5 0 4	コンフュージョンモデル	
5 0 6	言語モデル	50

- 5 0 8 辞書
- 5 1 0 単語誤り率計算機
- 6 0 0 状態
- 6 0 2 状態
- 6 0 4 状態
- 6 0 6 状態
- 6 0 8 状態
- 6 1 0 状態
- 6 1 2 最終状態

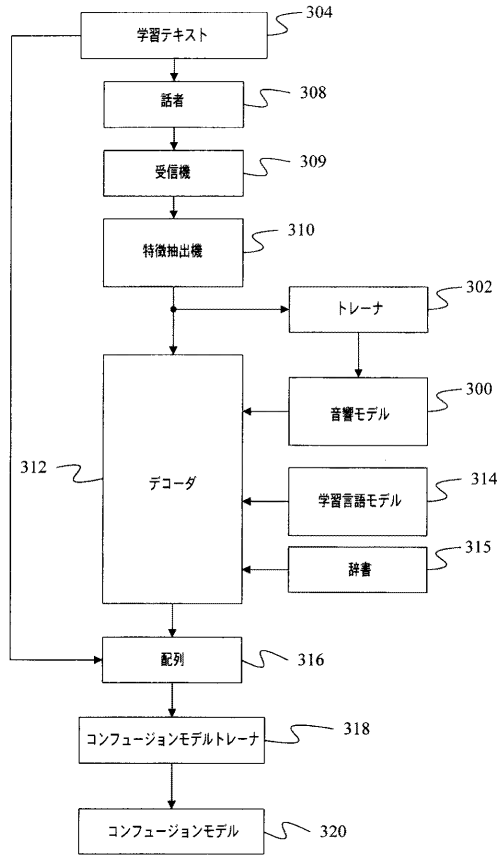
【図1】



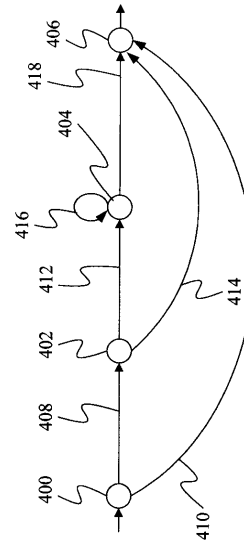
【図2】



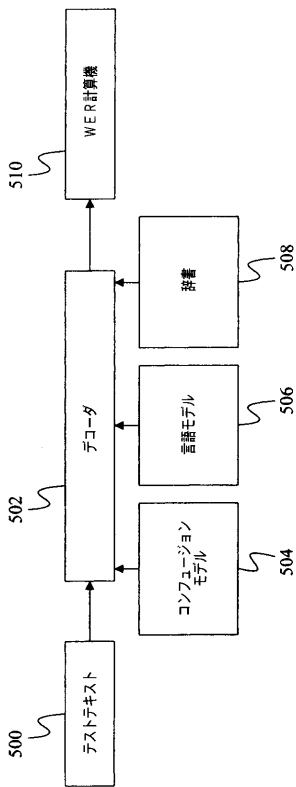
【図3】



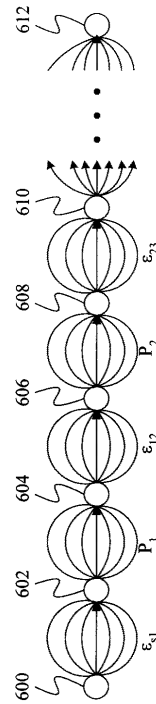
【図4】



【図5】



【図6】



フロントページの続き

- (72)発明者 デン ヤンガン
アメリカ合衆国 21204 メリーランド州 トーソン スティーブンソン レーン 340
アパートメント ナンバーシー4
- (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 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 チブリアン チェルバ
アメリカ合衆国 98112 ワシントン州 シアトル 43 アベニュー イースト 1928
アパートメント 4

審査官 井上 健一

- (56)参考文献 特表2003-535410(JP,A)
特開平06-318096(JP,A)
特開平10-097271(JP,A)
特開2004-069858(JP,A)

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