

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2004-118461
(P2004-118461A)

(43) 公開日 平成16年4月15日(2004.4.15)

(51) Int. Cl. ⁷	F I	テーマコード (参考)
G06F 17/22	G06F 17/22 514U	5B009
G06F 17/21	G06F 17/21 550A	

審査請求 未請求 請求項の数 10 O L (全 17 頁)

(21) 出願番号	特願2002-279934 (P2002-279934)	(71) 出願人	391055933 マイクロソフト コーポレイション MICROSOFT CORPORATION 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 (番地なし)
(22) 出願日	平成14年9月25日 (2002. 9. 25)	(74) 代理人	100077481 弁理士 谷 義一
		(74) 代理人	100088915 弁理士 阿部 和夫
		(74) 代理人	100106998 弁理士 橋本 博一

最終頁に続く

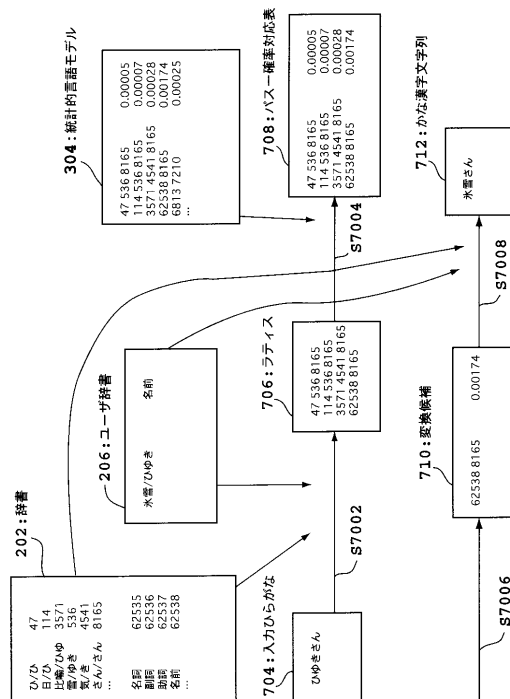
(54) 【発明の名称】 言語モデルのトレーニング方法、かな漢字変換方法、言語モデルのトレーニング装置、かな漢字変換装置、コンピュータプログラムおよびコンピュータ読み取り可能な記録媒体

(57) 【要約】

【課題】 統計的言語モデルに基づきながら、品詞によって定義された語も変換できるようにしたかな漢字変換装置を提供する。

【解決手段】 ステップS7002において、コンピュータシステムは、辞書202とユーザ辞書とを用いて、入力ひらがな704から語IDと品詞IDの混ざったIDの組み合わせ(パス)の集合(ラティス)を作成する。ステップS7004において、統計的言語モデル304から、各パスが生起する確率を取り出し、各パスと確率を対応付けたパス-確率対応表708を生成する。ステップS7006では、パス-確率対応表708のうちから、最も確率の高いパスを変換候補710として選択する。そして、ステップS7008では、辞書202とユーザ辞書206とを用いて、選択されたパスをかな漢字文字列712に変換する。

【選択図】 図7



【特許請求の範囲】

【請求項 1】

文字列を使用したコンピュータによる言語モデルのトレーニング方法であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、前記文字列に含まれる語および品詞に識別子を付与するステップと、前記文字列の中で、前記付与された識別子が特定の順序で生起する確率を示す言語モデルを生成するステップとを備えることを特徴とする言語モデルのトレーニング方法。

【請求項 2】

請求項 1 に記載の言語モデルのトレーニング方法により生成された言語モデルを用いたコンピュータによるかな漢字変換方法であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、入力された文字列から、識別子の組み合わせの集合を作成するステップと、前記言語モデルから、前記作成された識別子の組み合わせの各々が生起する確率を取り出すステップと、前記取り出された確率が最も高い識別子の組み合わせを選択するステップと、前記辞書を用いて、前記選択された識別子の組み合わせをかなまたは漢字に変換するステップとを備えることを特徴とするかな漢字変換方法。

【請求項 3】

請求項 2 に記載のかな漢字変換方法において、前記コンピュータは語と品詞との対応がユーザまたはベンダにより登録される登録辞書を有し、前記変換するステップは、前記辞書と前記登録辞書とを用いて前記識別子の組み合わせを変換することを特徴とするかな漢字変換方法。

【請求項 4】

語および品詞を含む文字列を使用した言語モデルのトレーニング装置であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を記憶した記憶手段と、前記記憶手段に記憶された辞書を用いて、前記文字列に含まれる語および品詞に識別子を付与する付与手段と、前記文字列の中で、前記付与手段により付与された識別子が特定の順序で生起する確率を示す言語モデルを生成する生成手段とを備えることを特徴とする言語モデルのトレーニング装置。

【請求項 5】

請求項 4 に記載の言語モデルのトレーニング装置により生成された言語モデルを用いたかな漢字変換装置であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を記憶した記憶手段と、前記記憶手段に記憶された辞書を用いて、入力された文字列から、識別子の組み合わせの集合を作成する作成手段と、前記言語モデルから、前記作成手段により作成された識別子の組み合わせの各々が生起する確率を取り出す取出手段と、前記取出手段により取り出された確率が最も高い識別子の組み合わせを選択する選択手段と、前記辞書を用いて、前記選択手段により選択された識別子の組み合わせをかなまたは漢字に変換する変換手段とを備えることを特徴とするかな漢字変換装置。

【請求項 6】

請求項 5 に記載のかな漢字変換装置において、前記記憶手段は語と品詞との対応がユーザまたはベンダにより登録される登録辞書を更に記憶し、前記変換手段は、前記記憶手段に記憶された辞書と登録辞書とを用いて前記識別子の組み合わせを変換することを特徴とす

10

20

30

40

50

るかな漢字変換装置。

【請求項 7】

語および品詞を含む文字列を使用した言語モデルのトレーニングを行うためのコンピュータプログラムであって、コンピュータに対し、
表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、前記文字列に含まれる語および品詞に識別子を付与するステップと、
前記文字列の中で、前記付与された識別子が特定の順序で生起する確率を示す言語モデルを生成するステップと
を実行させることを特徴とするコンピュータプログラム。

【請求項 8】

請求項 7 に記載のコンピュータプログラムにより生成されたされた言語モデルを有するコンピュータに対し、
表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、入力された文字列から、識別子の組み合わせの集合を作成するステップと、
前記言語モデルから、前記作成された識別子の組み合わせの各々が生起する確率を取り出すステップと、
前記取り出された確率が最も高い識別子の組み合わせを選択するステップと、
前記辞書を用いて、前記選択された識別子の組み合わせをかなまたは漢字に変換するステップと
を実行させることを特徴とするコンピュータプログラム。

【請求項 9】

請求項 8 に記載のコンピュータプログラムにおいて、前記コンピュータは語と品詞との対応がユーザまたはベンダにより登録される登録辞書を有し、前記変換するステップは、前記辞書と前記登録辞書とを用いて前記識別子の組み合わせを変換することを特徴とするコンピュータプログラム。

【請求項 10】

請求項 7 から 9 のいずれかに記載のコンピュータプログラムを記録したことを特徴とするコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、言語モデルのトレーニング方法、かな漢字変換方法、言語モデルのトレーニング装置、かな漢字変換装置、コンピュータプログラムおよびコンピュータ読み取り可能な記録媒体に関し、より詳細には、語と品詞の混ざった文字列を使用して統計的言語モデルで品詞を扱う言語モデルのトレーニング方法、かな漢字変換方法、言語モデルのトレーニング装置、かな漢字変換装置、コンピュータプログラムおよびコンピュータ読み取り可能な記録媒体に関する。

【0002】

【従来の技術】

日本語の文字列を入力する装置として、キーボードから入力したいかな漢字文字列に対応するかな文字列を入力し、漢字変換キーの入力にตอบสนองして、かな文字列をかな漢字文字列に変換するパーソナルコンピュータやワード・プロセッサなどのかな漢字変換装置が従来から知られている。この装置に入力したかな文字列をかな漢字文字列に変換する場合は、漢字変換用の特定の 1 つまたは複数のキーを組み合わせで押下し、かな漢字文字列の候補を表示する。また、連続して候補を表示することも可能であり、この場合前候補キーや次候補キーを押下するなどして、漢字を確定することができる。かな文字列をカタカナ文字列に変換する場合や、ローマ字文字列に変換する場合も、上記と同様の手順で行われる。

【0003】

入力された文字列についてかな漢字変換を行うかな漢字変換装置では、文字列に対応するかな漢字の候補を決定するのに、形態素や各フレームの解析情報を参照することによって

10

20

30

40

50

、変換精度を高めている。形態素とは、1つ以上の音素からなる意味をもった最小の言語単位をいい、形態素解析では、文字列に含まれている形態素の切れ目を認識し、および形態素の品詞を認定する。

【0004】

このような日本語の形態素解析において、従来から接続コスト最小法に基づく変換処理が知られている。これは、文の単語分割に対して何らかの接続コストを設定し、文全体で接続コストの和が最小となるような単語分割を選択する方法である。従って、表記、読み、品詞さえ指定されれば、どのような語も変換できる（例えば、非特許文献1参照）。

【0005】

接続コスト最小法では、品詞接続コストと単語コストを定義する。品詞接続コストは接続がまれな品詞間ほど大きく、単語コストは出現頻度が小さいほど大きくなるように設定する。しかし、接続コスト最小法では、コストを設定するための方法論が存在しない。

【0006】

この欠点を解消するための方法として、統計的言語モデルの研究が盛んに行われている（例えば、非特許文献1参照）。これは、接続コスト最小法のコストに相当する言語モデルを対象領域のテキストから自動的に学習する方法であり、情報理論と確率理論とに基づく明確な理論的根拠を備え、かつ実験的にも高い精度を持っている。

【0007】

【非特許文献1】

田中穂積監修、「自然言語処理 - 基礎と応用 - 」電子情報通信学会、平成11年3月25日

【0008】

【発明が解決しようとする課題】

統計的言語モデルに基づいたかな漢字変換システムでは、実世界において各語が生起する確率が指定されなければならない。しかしながら、各語が生起する確率を知りようがないユーザにとっては、辞書に対して確率を指定することは不可能である。このため、語に対して表記、読み、品詞が指定されていても、実世界における確率が指定されていなければ、変換処理を行うことができないという問題があった。

【0009】

本発明はこのような問題に鑑みてなされたものであり、その目的とするところは、統計的言語モデルに基づきながら、品詞によって定義された語も変換できるかな漢字変換装置およびかな漢字変換方法を提供することにある。

【0010】

【課題を解決するための手段】

このような目的を達成するために、請求項1に記載の発明は、文字列を使用したコンピュータによる言語モデルのトレーニング方法であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、前記文字列に含まれる語および品詞に識別子を付与するステップと、前記文字列の中で、前記付与された識別子が特定の順序で生起する確率を示す言語モデルを生成するステップとを備えることを特徴とする。

【0011】

また、請求項2に記載の発明は、請求項1に記載の言語モデルのトレーニング方法により生成された言語モデルを用いたコンピュータによるかな漢字変換方法であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、入力された文字列から、識別子の組み合わせの集合を作成するステップと、前記言語モデルから、前記作成された識別子の組み合わせの各々が生起する確率を取り出すステップと、前記取り出された確率が最も高い識別子の組み合わせを選択するステップと、前記辞書を用いて、前記選択された識別子の組み合わせをかなまたは漢字に変換するステップとを備えることを特徴とする。

【0012】

また、請求項3に記載の発明は、請求項2に記載のかな漢字変換方法において、前記コン

10

20

30

40

50

コンピュータは語と品詞との対応がユーザまたはベンダにより登録される登録辞書を有し、前記変換するステップは、前記辞書と前記登録辞書とを用いて前記識別子の組み合わせを変換することを特徴とする。

【0013】

また、請求項4に記載の発明は、語および品詞を含む文字列を使用した言語モデルのトレーニング装置であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を記憶した記憶手段と、前記記憶手段に記憶された辞書を用いて、前記文字列に含まれる語および品詞に識別子を付与する付与手段と、前記文字列の中で、前記付与手段により付与された識別子が特定の順序で生起する確率を示す言語モデルを生成する生成手段とを備えることを特徴とする。

10

【0014】

また、請求項5に記載の発明は、請求項4に記載の言語モデルのトレーニング装置により生成された言語モデルを用いたかな漢字変換装置であって、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を記憶した記憶手段と、前記記憶手段に記憶された辞書を用いて、入力された文字列から、識別子の組み合わせの集合を作成する作成手段と、前記言語モデルから、前記作成手段により作成された識別子の組み合わせの各々が生起する確率を取り出す取出手段と、前記取出手段により取り出された確率が最も高い識別子の組み合わせを選択する選択手段と、前記辞書を用いて、前記選択手段により選択された識別子の組み合わせをかなまたは漢字に変換する変換手段とを備えることを特徴とする。

20

【0015】

また、請求項6に記載の発明は、請求項5に記載のかな漢字変換装置において、前記記憶手段は語と品詞との対応がユーザまたはベンダにより登録される登録辞書を更に記憶し、前記変換手段は、前記記憶手段に記憶された辞書と登録辞書とを用いて前記識別子の組み合わせを変換することを特徴とする。

【0016】

また、請求項7に記載の発明は、語および品詞を含む文字列を使用した言語モデルのトレーニングを行うためのコンピュータプログラムであって、コンピュータに対し、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、前記文字列に含まれる語および品詞に識別子を付与するステップと、前記文字列の中で、前記付与された識別子が特定の順序で生起する確率を示す言語モデルを生成するステップとを実行させることを特徴とする。

30

【0017】

また、請求項8に記載の発明は、請求項7に記載のコンピュータプログラムにより生成されたされた言語モデルを有するコンピュータに対し、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、入力された文字列から、識別子の組み合わせの集合を作成するステップと、前記言語モデルから、前記作成された識別子の組み合わせの各々が生起する確率を取り出すステップと、前記取り出された確率が最も高い識別子の組み合わせを選択するステップと、前記辞書を用いて、前記選択された識別子の組み合わせをかなまたは漢字に変換するステップとを実行させることを特徴とする。

40

【0018】

また、請求項9に記載の発明は、請求項8に記載のコンピュータプログラムにおいて、前記コンピュータは語と品詞との対応がユーザまたはベンダにより登録される登録辞書を有し、前記変換するステップは、前記辞書と前記登録辞書とを用いて前記識別子の組み合わせを変換することを特徴とする。

【0019】

また、請求項10に記載の発明は、コンピュータ読み取り可能な記録媒体であって、請求項7から9のいずれかに記載のコンピュータプログラムを記録したことを特徴とする。

【0020】

【発明の実施の形態】

50

本発明の好ましい実施形態を示す以下の説明では、本明細書の一部を形成し、本発明を実践できる特定の実施形態を示す添付図面を参照する。本発明の範囲から逸脱せずに、他の実施形態を使用することができ、構造的変更を行うことができることを理解されたい。

【0021】

図1に、本発明を実施することができる適切なコンピュータシステム100の一例を示す。このコンピュータシステム100は適切なコンピュータシステムの一例にすぎず、本発明の使用法または機能性の範囲に関するいかなる限定をも示唆しようとするものではない。また、コンピュータシステム100は、同図に示す任意の1つまたは複数の構成要素の組み合わせに関する依存または要件を有するものと解釈されるべきではない。

【0022】

本発明は、コンピュータで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なコンテキストに関して説明することができる。一般に、プログラムモジュールは、特定タスクを実行するか、または特定の抽象データタイプを実施するルーチン、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、通信ネットワークを介してリンクされている遠隔処理デバイスによってタスクが実施される分散型コンピュータシステムで実施することもできる。分散型コンピュータシステムでは、プログラムモジュールは、記憶装置を含めて、ローカルおよび遠隔コンピュータの記録媒体内に配置することができる。

10

【0023】

図1を参照すると、本発明を実施するための例示的システムは、コンピュータ110の形態による汎用コンピューティングデバイスを含む。コンピュータ110の構成要素は、限定はしないが、処理ユニット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ユニット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いずれかを使用するメモリバスまたはメモリコントローラ、周辺バス、およびローカルバスを含む複数バス構造のうちのどれであっても良い。限定的ではないが、例示として、このようなアーキテクチャには、Industry Standard Architecture (ISA)バス、Micro Channel Architecture (MCA)バス、Enhanced ISA (EISA)バス、Video Electronics Standards Associate (VESA)ローカルバス、およびMezzanineバスとしても知られているPeripheral Component Interconnect (PCI)バスが含まれる。

20

30

【0024】

コンピュータ110は、一般に様々なコンピュータ読み取り可能な記録媒体を含む。コンピュータ読み取り可能な記録媒体は、コンピュータ110によってアクセス可能ないかなる使用可能な媒体であってもよく、揮発性媒体および不揮発性媒体、取り外し可能媒体および取り外し不可能媒体の両方を含むものである。限定はしないが例示として、コンピュータ読み取り可能な記録媒体は、コンピュータ読み取り可能な記録媒体と通信媒体とを含むことができる。

【0025】

コンピュータ読み取り可能な記録媒体は、コンピュータ実行可能命令、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格納用のいかなる方法または技術で実施される揮発性および不揮発性、取り外し可能および取り外し不可能媒体の両方を含む。コンピュータ読み取り可能な記録媒体は、限定はしないが、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、デジタル多目的ディスク(DVD)または他の光ディスクストレージ、磁気カセット、磁気テープ、磁気ディスクストレージまたは他の磁気記憶装置、所望の情報を格納するために使用され、またコンピュータ110によってアクセスすることができる他のいかなる媒体をも含む。

40

【0026】

通信媒体は、通常、コンピュータ実行可能命令、データ構造、プログラムモジュールまた

50

は搬送波または他の搬送メカニズムなどの変調されたデータ信号形式の他のデータを組み込むものであり、いかなる情報伝達媒体をも含むものである。「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的な設定を有する信号、または信号中に情報を符号化するような方法で変更された信号を意味している。限定はしないが例示として、通信媒体は、有線ネットワークまたは直接有線接続などの有線媒体と、音波、RF、赤外線または他の無線媒体などの無線媒体を含む。コンピュータ読み取り可能な記録媒体の範囲には、上記のいかなる組み合わせも含まれるべきである。

【0027】

システムメモリ130は、読み出し専用メモリ(ROM)131およびランダムアクセスメモリ(RAM)132などの揮発性または不揮発性メモリ形式のコンピュータ読み取り可能な記録媒体を含む。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ために役立つ基本ルーチンを含んでいる基本入出力システム133(BIOS)は、一般にROM131内に格納されている。RAM132は、一般に、処理ユニット120に即時アクセス可能な、またはその時点において処理ユニット120により操作されているデータまたはプログラムモジュールを含む。限定はしないが例示として、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

【0028】

コンピュータ110は、他の取り外し可能/取り外し不可能、揮発性/不揮発性コンピュータ読み取り可能な記録媒体を含むこともできる。例示として、図1に、取り外し不可能な不揮発性磁気媒体に対して読み出しまたは書込みするハードディスクドライブ141、取り外し可能な不揮発性磁気ディスク152に対して読み出しまたは書込みする磁気ディスクドライブ151、取り外し可能な不揮発性光ディスク156に対して読み出しまたは書込みをする、CD-ROMまたは他の光学媒体などの光ディスクドライブ155を示す。例示的オペレーティング環境で使用することができる他の取り外し可能/取り外し不可能な揮発性/不揮発性コンピュータ読み取り可能な記録媒体は、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、デジタル多目的ディスク、デジタルビデオテープ、半導体RAM、半導体ROMなどを含む。ハードディスクドライブ141は、一般に、インターフェース140などの取り外し不可能なメモリインターフェースを介してシステムバス121に接続されており、磁気ディスクドライブ151と光ディスクドライブ155は、一般に、インターフェース150などの取り外し可能なメモリインターフェースによってシステムバス121に接続されている。

20

30

【0029】

図1に示す上記のドライブおよびそれらに関連したコンピュータ読み取り可能な記録媒体は、コンピュータ実行可能命令、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のための他のデータの記憶装置を提供する。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を格納してい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良いという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、その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、最低限、それらが異なるコピーであることを示すために異なる番号が与えられている。ユーザは、キーボード162および、一般にマウス、トラックボールまたはタッチパッドと呼ばれるポインティングデバイス161などの入力デバイスによってコンピュータ110にコマンドおよび情報を入力することができる。他の入力デバイス(図示せず)は、マイクロフォン、ジョイスティック、ゲームパッド、衛星放送アンテナ、スキャナなどを含むことができる。これらのおよび他の入力デバイスは、システムバスに結合されたユーザ入力インターフェース160を介して処理ユニット120に接続されることがしばしばあるが、パラレルポート、ゲームポートまたはユニ

40

50

バーサルシリアルバス（USB）などの他のインターフェースおよびバス構造に接続されることもできる。モニタ191または他のタイプの表示装置も、ビデオインターフェース190などのインターフェースを介してシステムバス121に接続される。モニタの他に、コンピュータは、出力周辺インターフェース195を介して接続することができるスピーカ197およびプリンタ196など、他の周辺出力装置を含むこともできる。

【0030】

コンピュータ110は、遠隔コンピュータ180などの1つまたは複数の遠隔コンピュータへの論理接続を使用してネットワーク化された環境において動作することができる。遠隔コンピュータ180は、別のパーソナルコンピュータ、サーバ、ルータ、ネットワークPC、ピアデバイスまたは他の共通ネットワークノードであってよく、図1では記憶装置181しか示していないが、一般にコンピュータ110に関して上記で説明した多くのまたはすべての要素を含む。図1で示す論理接続は、ローカルエリアネットワーク（LAN）171およびワイド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むこともできる。このようなネットワーキング環境は、事務所、企業全体に巡らされている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一般的なことである。

10

【0031】

LANネットワーキング環境で使用されるとき、コンピュータ110は、ネットワークインターフェースまたはアダプタ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されるとき、コンピュータ110は、通常、インターネットなどのWAN173を介して通信を確立するモデム172または他の手段を含む。内部であっても外部であっても良いモデム172は、ユーザ入力インターフェース160または他の適切なメカニズムを介してシステムバス121に接続することができる。ネットワーク環境において、コンピュータ110に関して示されたプログラムモジュールまたはその一部は、遠隔記憶装置に格納することができる。限定はしないが例示として、図1に、記憶装置181上に常駐するものとして遠隔アプリケーションプログラム185を示す。図示するネットワーク接続は一例であり、コンピュータ間で通信リンクを確立する他の手段を使用することもできることを理解されたい。

20

【0032】

以下の説明では、本発明は、特に指摘しない限り、アプリケーションプログラムのコンピュータ実行可能命令をシステムメモリ131にロードした処理ユニット120が、そのコンピュータ実行可能命令に基づき実行することができる動作を説明する。この動作において、処理ユニット120はコンピュータ実行可能命令に基づきプログラムデータ137を参照し、あるいはその更新を行う。

30

【0033】

従って、コンピュータによって実行されるときに時折表現されることのあるこのような動作および演算には、コンピュータの処理ユニットによる、構造化形式のデータを表現する電子信号の操作が含まれることを理解されたい。この操作はデータを変換するか、またはコンピュータのメモリシステム中の記憶場所にデータを維持し、そこで、当業者が良く理解している方法でコンピュータの演算を再構成あるいは変更する。データが維持されているデータ構造は、データ形式によって規定される特定の特性を有するメモリの物理的な記憶場所である。本発明を上記の条件で説明してはいるが、以下で説明する様々な動作および演算はハードウェアでも実施可能であることを当業者なら理解するように、この説明は限定を意図するものではない。

40

【0034】

図2は、本実施形態に係るプログラムデータ137の内容をより詳細に示す図であり、本発明に関わる部分のみを概略的に示している。

【0035】

プログラムデータ137は、コーパス202と、辞書204と、ユーザ辞書206とを含んでいる。コーパス202は、自然言語処理等に利用される大規模テキストデータであっ

50

て、文字列が形態素ごとに分割され、各形態素について品詞が決定された（即ち、品詞タグ付けされた）ものである。その他、係り受けなどの統語情報が付加されたものもコーパス202として利用することができる。辞書204は、語および品詞の各々に対する識別子（ID）を定義したデータである。ここで、語は文字の表記とその読みとを含んでいる。ユーザ辞書206は、ユーザ個人が使い勝手を良くするために単語や定型句を登録して作成する登録辞書の1つである。ここで、登録辞書は、ユーザ辞書の他、専門辞書や分野別辞書などのベンダにより登録されるものであっても良い。

【0036】

図3は、本実施形態に係るコンピュータシステムによる言語モデルのトレーニング方法の概要を示す図である。以下では説明を簡単にするためにバイグラムモデルを例に挙げて説明するが、本発明はトライグラムモデル等の他のマルチグラムモデルにも適用可能であることはいうまでもない。

10

【0037】

まず、コーパス202に含まれている文字列から、表記と読みとを含む語のID、および、品詞のIDを定義した辞書204を用いることにより、特定のIDペアの生起回数を示すIDペア - 生起回数対応表302を生成する（S3002）。このIDペアの生起回数から、各IDペアがコーパス202内で生起する確率を示す統計的言語モデル304が生成される（S3004）。以下、図4～図6を参照し、上述した方法の各ステップについて詳細に説明する。

【0038】

図4に示すように、辞書204は、語データ402と、品詞データ404とを有する。語データ402は、「は/は」、「だ/だ」、「今日/きょう」、のように、表記と読みとを含む語406の集合である語データ402と、ID408との対応を示すものである。また、語データ402には、文頭および文末とID408との対応も定義されている。品詞データ404は、名詞、副詞、助詞といった品詞410の各々のID412を示すものである。

20

【0039】

いま、図5に示すようにコーパス202が
文頭 今日/きょう/名詞 は/は/助詞 天気/てんき/副詞 だ/だ/助動詞 。 /
。 /句点 文末

30

という情報を含んでいるものとする。ステップS3002において、コンピュータシステムは辞書204を使用し、コーパス202に含まれる文字列「今日は天気だ。」に含まれる語および品詞に、辞書204内のIDを付与する。次いで、コーパス202内で、特定のIDペアが生起する回数510を数える。ここで、特定のIDペアとは、語IDと語IDのペア502、品詞IDと品詞IDのペア504、語IDと品詞IDのペア506、および品詞IDと語IDのペア508である。

【0040】

次いで、ステップS3004では、図6に示すように、IDペア - 生起回数対応表302内の各IDペアの生起回数510に基づき、特定のIDのペアが生起する確率602を示す統計的言語モデル304を生成する。

40

【0041】

以上説明したように、本実施形態では、コーパス内の語と品詞とが混ざった少なくとも1つの文字列から、語と品詞のペアが生起する確率、品詞同士のペアが生起する確率、および語同士のペアが生起する確率を計算する。このように言語モデルをトレーニングすることにより、語の確率情報による変換とともに品詞情報を用いた変換も可能となる。以下、上述したように生成された言語モデルを用いたかな漢字変換方法について説明する。

【0042】

図7は、本実施形態に係るかな漢字変換方法の概要を示す図である。なお、以下の説明において使用されるコンピュータシステムとして図1に示すものが使用されるが、これは上述した言語モデルのトレーニング方法に使用されるものと同様であっても良く、異なるも

50

のであっても良い。後者の場合、上述のように生成された言語モデルは、CD-ROM等の取り外し可能揮発性メモリ、あるいは有線ネットワーク等の通信媒体を使用して、以下の処理を実行するコンピュータシステムへ提供することができる。

【0043】

ステップS7002において、コンピュータシステムは、辞書202とユーザ辞書とを用いて、入力ひらがな704から語IDと品詞IDの混ざったIDの組み合わせ(パス)の集合(ラティス)を作成する。ステップS7004において、統計的言語モデル304から、各パスの生起する確率を取り出し、各パスと確率を対応付けたパス-確率対応表708を生成する。ステップS7006では、パス-確率対応表708のうちから、最も確率の高いパスを変換候補710として選択する。そして、ステップS7008では、辞書202とユーザ辞書206とを用いて、選択されたパスをかな漢字文字列712に変換する。以下、図8~10を参照し、本実施形態に係るかな漢字変換方法について詳細に説明する。

10

【0044】

まず、図8に示すように、キーボード162等の入力手段を介して、品詞と語とを含む入力ひらがな704がコンピュータシステムに入力される。入力ひらがな704は、入力と同時にモニタ191に表示される。ステップS7002では、入力ひらがな704から各種の形態素が生成され、辞書202およびユーザ辞書206からあらゆる可能性の変換すべき語の候補が取り出され、語IDまたは品詞IDの組み合わせ(パス)の集合であるラティス706を生成する。

20

【0045】

例えば、1文字目「ひ」の候補は「ひ/ひ/47」と「日/ひ/114」の2つである。また、これに続く候補は「雪/ゆき/536」である。最後の2文字の候補は「さん/さん/8165」である。ここまでで、

47 536 8165

114 536 8165

というパスができる。また、同時にユーザ辞書やベンダ辞書も参照され、IDで候補が作成され、ラティスに追加される。同図に示す例では、ユーザ辞書206中に、「ひゆき」という読みに対して「氷雪/名前」という表記および品詞が対応付けられている。辞書202によれば「名前」という品詞のIDは62538なので、このIDと「さん/さん/8165」という語のIDと組み合わせた

30

62538 8165

というパスがラティス706に追加される。なお、辞書202およびユーザ辞書206を参照する順序は逆であっても良い。

【0046】

ステップS7004では、図9に示すように、統計的言語モデル304からラティス706に含まれるパスの生起する確率を取り出される。例えば、最初のパスは

47 536 8165

なので、その確率0.00005が取り出される。同様にして、ラティス706を構成する各パスの確率を取り出され、パス-確率対応表708が生成される。

40

【0047】

ステップS7006では、図10に示すように、パス-確率対応表708のうちから、最も確率が高いパス

62538 8165

を変換候補710として選択する。

【0048】

そして、図11に示すように、ステップS7008でまず変換候補710の最初のIDである62538に対応する語または品詞を辞書202から取り出す。辞書202を参照すると、62538は「名前」という品詞であることが分かる。本実施形態において、「名前」はユーザ辞書206に登録される品詞であることから、次にユーザ辞書206が参照

50

され、62538というIDは、入力文字列「ひゆき」に対応する表記「氷雪」に変換される。

【0049】

次のIDは8165なので、このIDに対応する語または品詞が辞書202から取り出される。ここで、8165というIDは「さん」という表記に変換される。このようにして、最終的に「氷雪さん」という文字列が得られ、モニタ191にかな漢字文字列712が表示される。

【0050】

以上、本発明の好適な実施の形態について説明したが、本発明の前述の説明は、例示および説明を目的として提示されたものである。網羅的であること、または本発明を開示された正確な形態に制限することは、意図されていない。多数の修正形態および変形形態が、上の教示に鑑みて可能である。本発明の範囲は、この詳細な説明によるのではなく、請求項によって制限されることが意図されている。

10

【0051】

【発明の効果】

以上説明したように、本発明では、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、文字列に含まれる語および品詞に識別子を付与するステップと、文字列の中で、付与された識別子が特定の順序で生起する確率を示す言語モデルを生成するステップとを備える。また、表記と読みとを含む語の識別子、および、品詞の識別子を定義した辞書を用いて、入力された文字列から、識別子の組み合わせの集合を作成するステップと、言語モデルから、作成された識別子の組み合わせの各々が生起する確率を取り出すステップと、取り出された確率が最も高い識別子の組み合わせを選択するステップと、辞書を用いて、選択された識別子の組み合わせをかなまたは漢字に変換するステップとを備える。従って、語の確率情報による変換とともに品詞情報を用いた変換も可能となる。

20

【0052】

これは、特にサイズの小さい統計的言語モデルに有効である。統計的言語モデルのサイズを小さくするということは、即ち変換精度を落とすことを意味する。そこで、本発明を適用することにより、基本語ではない語を品詞IDを使用して変換できるので、変換精度を高く保つことができる。

30

【0053】

また、コンピュータは語と品詞との対応がユーザまたはベンダにより登録される登録辞書を有し、変換するステップは、辞書と登録辞書とを用いて識別子の組み合わせを変換するので、統計的言語モデルにおいて語IDだけでなく品詞IDもラティスに追加することで、ユーザ辞書やベンダ辞書に含まれている語の変換が可能となるという効果を奏する。

【図面の簡単な説明】

【図1】本発明を実施する例示的システムを構成するコンピュータシステムを示す図である。

【図2】本発明実施形態によるプログラムデータの内容を示すブロック図である。

【図3】本発明実施形態による言語モデルのトレーニング方法の動作の概要を示す図である。

40

【図4】本発明実施形態による言語モデルのトレーニング方法の動作を示す図である。

【図5】本発明実施形態による言語モデルのトレーニング方法の動作を示す図である。

【図6】本発明実施形態による言語モデルのトレーニング方法の動作を示す図である。

【図7】本発明実施形態によるかな漢字変換方法の動作の概要を示す図である。

【図8】本発明実施形態によるかな漢字変換方法の動作を示す図である。

【図9】本発明実施形態によるかな漢字変換方法の動作を示す図である。

【図10】本発明実施形態によるかな漢字変換方法の動作を示す図である。

【図11】本発明実施形態によるかな漢字変換方法の動作を示す図である。

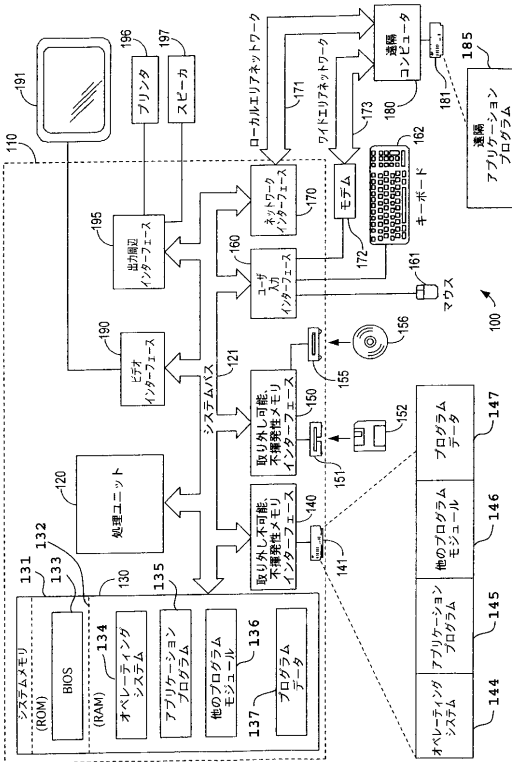
【符号の説明】

50

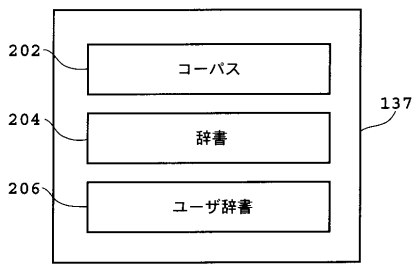
1 0 0	コンピュータシステム	
1 1 0	コンピュータ	
1 2 0	処理ユニット	
1 2 1	システムバス	
1 3 0	システムメモリ	
1 3 1	読み出し専用メモリ	
1 3 2	ランダムアクセスメモリ	
1 3 3	基本入出力システム	
1 3 4	オペレーティングモジュール	
1 3 5	アプリケーションプログラム	10
1 3 6	他のプログラムモジュール	
1 3 7	プログラムデータ	
1 4 0	取り外し不可能不揮発性メモリインターフェース	
1 4 1	ハードディスクドライブ	
1 4 4	オペレーティングシステム	
1 4 5	アプリケーションプログラム	
1 4 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	取り外し可能不揮発性メモリインターフェース	
1 5 1	磁気ディスクドライブ	20
1 5 2	取り外し可能な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	取り外し可能な不揮発性光ディスク	
1 6 0	ユーザ入力インターフェース	
1 6 1	ポインティングデバイス	
1 6 2	キーボード	
1 7 0	アダプタ	
1 7 1	ローカルエリアネットワーク (LAN)	
1 7 2	モデム	
1 7 3	ワイドエリアネットワーク (WAN)	30
1 8 0	遠隔コンピュータ	
1 8 1	記憶装置	
1 8 4	マルチレベルキャッシュ	
1 8 5	遠隔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ース	
1 9 1	モニタ	
1 9 5	出力周辺インターフェース	
1 9 6	プリンタ	
1 9 7	スピーカ	
2 0 2	コーパス	40
2 0 4	辞書	
2 0 6	ユーザ辞書	
3 0 2	IDペア - 生起回数対応表	
3 0 4	統計的言語モデル	
4 0 2	語データ	
4 0 4	品詞データ	
4 0 6	表記 / 読み	
4 0 8	ID番号	
4 1 0	品詞	
4 1 2	ID番号	50

- 5 0 2 語 I D - 語 I D の ペア
- 5 0 4 品 詞 I D - 品 詞 I D の ペア
- 5 0 6 語 I D - 品 詞 I D の ペア
- 5 0 8 品 詞 I D - 語 I D の ペア
- 5 1 0 I D ペアの生起回数
- 6 0 2 I D ペアの確率
- 7 0 4 入力ひらがな
- 7 0 6 ラティス
- 7 0 8 パス - 確率対応表
- 7 1 0 変換候補
- 7 1 2 かな漢字文字列

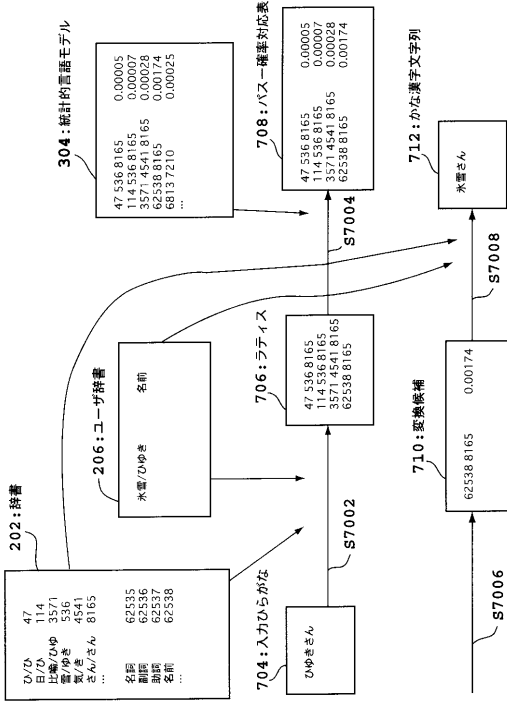
【 図 1 】



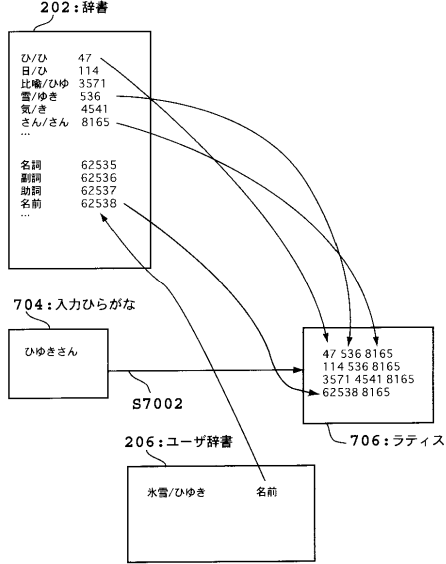
【 図 2 】



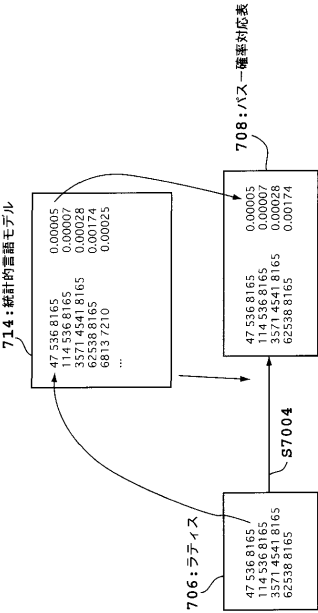
【 図 7 】



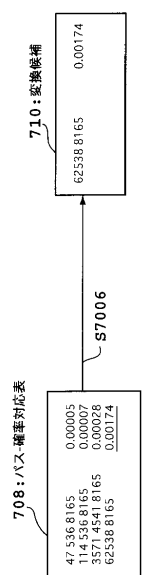
【 図 8 】



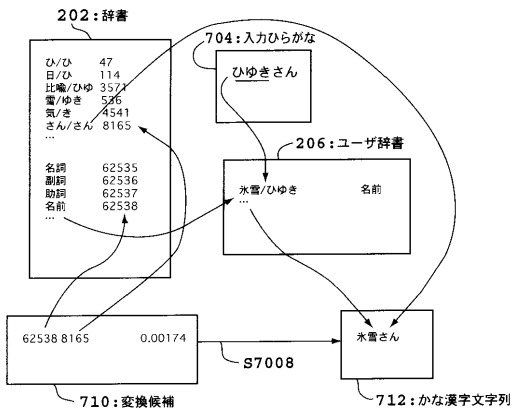
【 図 9 】



【 図 10 】



【図 11】



フロントページの続き

(72)発明者 石橋 紀子

東京都調布市調布ヶ丘 1 - 1 8 - 1 マイクロソフト株式会社 マイクロソフト調布技術センター
内

(72)発明者 鹿子木 宏明

東京都調布市調布ヶ丘 1 - 1 8 - 1 マイクロソフト株式会社 マイクロソフト調布技術センター
内

Fターム(参考) 5B009 MA05 QA01 TA09