

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4724357号
(P4724357)

(45) 発行日 平成23年7月13日(2011.7.13)

(24) 登録日 平成23年4月15日(2011.4.15)

(51) Int.Cl.

F I

G 0 6 F 17/28 (2006.01)

G 0 6 F 17/28

C

請求項の数 7 (全 17 頁)

(21) 出願番号	特願2003-125929 (P2003-125929)	(73) 特許権者	500046438
(22) 出願日	平成15年4月30日 (2003.4.30)		マイクロソフト コーポレーション
(65) 公開番号	特開2003-345796 (P2003-345796A)		アメリカ合衆国 ワシントン州 9805
(43) 公開日	平成15年12月5日 (2003.12.5)		2-6399 レッドモンド ワン マイ
審査請求日	平成18年4月26日 (2006.4.26)		クロソフト ウェイ
審査番号	不服2008-6010 (P2008-6010/J1)	(74) 代理人	100077481
審査請求日	平成20年3月10日 (2008.3.10)		弁理士 谷 義一
(31) 優先権主張番号	10/137,456	(74) 代理人	100088915
(32) 優先日	平成14年4月30日 (2002.4.30)		弁理士 阿部 和夫
(33) 優先権主張国	米国 (US)	(74) 復代理人	100115624
			弁理士 濱中 淳宏
		(74) 復代理人	100084191
			弁理士 合田 潔

最終頁に続く

(54) 【発明の名称】 コンピュータ可読媒体及び単語情報を得るコンピュータ実行方法並びに単語情報を格納する方法

(57) 【特許請求の範囲】

【請求項 1】

テキストアナライザとして動作するコンピュータにより実行され、複数の異なる自然言語処理において使用される複数のレキシコンからまとめて単語情報を取得するコンピュータ実行方法であって、

前記コンピュータは、

コンピュータ記憶媒体と、

プロセッサと、

前記コンピュータ記憶媒体に格納され、前記プロセッサ上で実行可能なプログラムとを備え、

前記プログラムは、複数の異なる自然言語処理を実行するために、前記プロセッサによりアクセス可能なコンピュータ記憶媒体に格納されたレキシコンから単語情報を取得する命令を含み、

各レキシコンは、

複数の単語を格納する単語リストセクションと、

複数の組のデータセクションであって、各組のデータセクションは、前記単語リストセクションの各単語に対応し、各組のデータセクションの各データセクションは、前記各単語について選択された実質的に異なる情報を格納するデータセクションと、

前記複数の組のデータセクションと分離した、前記単語リストセクションの各単語についての複数のポインタを格納する索引セクションであって、各複数のポインタは、ある自

然言語処理に関連付けられた第 1 の組のポインタと、異なる自然言語処理に関連付けられた第 2 の組のポインタとを有し、前記第 1 の組のポインタは、前記第 2 の組のポインタと異なり、各ポインタは、前記複数の組のデータセクションのデータをポイントする索引セクションとを備え、

前記プロセッサは、前記コンピュータ記憶媒体から前記プログラムを読み出し、前記プログラムを実行し、

前記方法は、

前記プロセッサが、類似した情報を有する前記各レキシコンの前記複数の組のデータセクションに選択的にアクセスするステップと、

前記プロセッサが、実行される特定の自然言語処理に応じて前記第 1 または第 2 の組のポインタを使用し、アクセスしたデータセクションから情報を取得するステップと
を含むことを特徴とするコンピュータ実行方法。

【請求項 2】

前記選択的にアクセスするステップは、類似する情報を有する前記各レキシコンの少なくとも 2 つのデータセクションの単語情報を組み合わせるステップを含むことを特徴とする請求項 1 に記載のコンピュータ実行方法。

【請求項 3】

前記選択的にアクセスするステップは、

第 1 のレキシコンのデータセクションから単語情報を取得するステップと、

第 2 のレキシコンのデータセクションから単語情報を取得するステップと、

前記第 2 のレキシコンのデータセクションの情報だけを使用するステップと
を含むことを特徴とする請求項 1 に記載のコンピュータ実行方法。

【請求項 4】

前記選択的にアクセスするステップは、停止インディケータが見つかるまで前記各レキシコンの類似するデータセクションから単語情報を取得するステップを含むことを特徴とする請求項 1 に記載のコンピュータ実行方法。

【請求項 5】

前記データセクションに選択的にアクセスするステップは、選択された順序で前記複数のレキシコンに順次アクセスするステップを含み、

前記プログラムは、各レキシコンからの単語情報を読み出すか否かの第 1 の命令を備え

、
前記プログラムは、前記複数のレキシコンの 2 つ以上からの単語情報を組み合わせる第 2 の命令を備えることを特徴とする請求項 1 に記載のコンピュータ実行方法。

【請求項 6】

前記データセクションに選択的にアクセスするステップは、

アクセスする各レキシコンについて、

特定の自然言語処理に応じて、対応する索引セクションのポインタ識別を確定するために、所与の単語に応じて対応する単語リストセクションにアクセスするステップと、

前記対応する索引セクションのポインタを得るために、前記ポインタ識別を使用するステップと、

前記複数のデータセクションのうちのどの対応するデータセクションに前記単語についての情報があるか、及び前記情報が前記対応するデータセクションのどこに位置するかを確定するために、前記ポインタを使用するステップと
を含むことを特徴とする請求項 1 に記載のコンピュータ実行方法。

【請求項 7】

前記選択的にアクセスするステップは、選択された順序で前記レキシコンに順次アクセスするステップを含むことを特徴とする請求項 1 に記載のコンピュータ実行方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、コンピュータ可読媒体及び単語情報を得るコンピュータ実行方法並びに単語情報を格納する方法に関する。より詳細には、言語またはテキストの処理に関連し、レキシコン (lexicon) を格納するための改良されたデータ構造、及びそのデータ構造を使用する方法に係る、レキシコンを有するコンピュータ可読媒体、単語情報を得るコンピュータ実行方法、単語情報を格納する方法及びそれらの方法を実施する命令を有するコンピュータ可読媒体に関する。

【 0 0 0 2 】

【従来の技術】

言語処理またはテキスト処理には多くのタイプのシステムが含まれる。例えば、パーサ、スペルチェッカ、文法チェッカ、ワードブレーカ、自然言語のプロセッサまたは理解システム、機械翻訳システムは、この広い範疇に該当するシステムのタイプのいくつかに過ぎない。

10

【 0 0 0 3 】

多くの言語またはテキスト処理システムに共通する重要なコンポーネントがレキシコンである。一般に、レキシコンは単語についての情報を含むデータ構造である。例えば、レキシコンは構文的情報及び意味的情報の指示を格納することができる。この例には、その単語が名詞であるか、動詞、形容詞であるかなどが挙げられる。また、異なるタイプの言語情報もレキシコンに格納することができる。しばしば、構文解析の助けとなる単語についての情報を格納するなど、特定タイプの言語処理に有用な他の情報を格納しておくことも有用である。さらに他のレキシコンでは、その単語が固有名詞か、地理的な場所かなどに

20

【 0 0 0 4 】

動作の際は、単語の入力文字列を受け取ると、言語またはテキスト処理システムはレキシコンにアクセスして、各単語についての格納された情報を得る。入力文字列中の各単語についての情報を集めると、言語またはテキスト処理システムはその入力文字列を処理するが、これには存在する可能性のあるあいまい性をその単語の情報に基づいて解消することが含まれる場合がある。例えば、自然言語処理システムでは、レキシコンは入力文字列中の各単語に品詞を割り当てる。次いで構文パーサがどの品詞の割り当てが適切であるかを判断し、入力文字列から構造を構築し、それを解釈のために意味コンポーネントに渡すことができる。

30

【 0 0 0 5 】

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)。

【 0 0 0 6 】

【非特許文献 1】

P.T.Sato 著「A COMMON PARSING SCHEME FOR LEFG-AND RIGHT-BRANCHING LANGUAGES」Computational Linguistics、Vol1.14、No.1、1988 年冬、p. 20 - 30

【 0 0 0 7 】

【発明が解決しようとする課題】

一般に、レキシコンの各項目は単一の大きなバイナリオブジェクトからなる。このフォーマットでは、情報へのアクセスは可能であるが、項目全体を読み込まずに、一般に使用される語彙情報への局所的なアクセスを容易には行うことができない。ある単語項目に関連するすべての情報をレキシコンから読み込まなければならない場合、特にその単語項目の情報のうちわずかな部分だけを必要とする場合には、より多くのメモリと処理時間が必要とされる。

40

【 0 0 0 8 】

語彙情報の変更または追加も難しい。具体的には、レキシコンを修正する、あるいはさらに情報を追加するために、レキシコンの作成者は、非常に複雑なデータ構造の整合性と編成を保ちながら、各項目中のすべてのビット、属性またはその他の情報を複製してから、所望の情報を変更するか、あるいは情報を追加しなければならない。

50

【 0 0 0 9 】

このため、上述の不都合点の１つ、一部、あるいはすべてに対処する改良したレキシコンデータ構造が必要とされる。

【 0 0 1 0 】

本発明は、このような課題に鑑みてなされたもので、その目的とするところは、レキシコンを格納するための高い柔軟性と効率を有する改良されたデータ構造を提供しそのデータ構造を使用可能とする、コンピュータ可読媒体及び単語情報を得るコンピュータ実行方法並びに単語情報を格納する方法を提供することにある。

【 0 0 1 1 】

【課題を解決するための手段】

本発明の一態様は、言語処理システムで使用するために適合された単語情報を有する、コンピュータ可読媒体に格納された単語レキシコンである。このレキシコンは、複数の単語を記憶する単語リストセクションと、それら複数の単語の単語情報を記憶する複数のデータセクションとを含む。複数のデータセクションは互いと単語リストセクションとから分離されている。単語情報にアクセスするために索引セクションが提供され、このセクションには複数のデータセクション中のデータをポイントするポインタが格納される。使用するポインタの識別は、単語リストセクション中の対応する単語に応じて決まる。

【 0 0 1 2 】

この改良されたレキシコン構造は、これまで得ることのできなかった柔軟性と効率を提供する。索引セクションと複数のデータセクションは、利用可能なコンピュータリソースなど言語処理システムの必要性に合わせて、レキシコンを適合することを可能にする。さらなる実施形態では、このレキシコン構造は、分類に基づいて単語情報を分類またはグループ化することを可能にする。例えば、この分類は、その単語項目が名詞か、動詞か、形容詞かなど、単語項目の品詞に基づくことができる。これにより、分類に応じて単語情報に選択的にアクセスすることができる。例示的实施形態では、対応する単語情報の分類を示す表示がポインタ中に提供される。

【 0 0 1 3 】

本発明の他の態様は、複数のデータセクションに単語情報を格納し、索引セクションにポインタ情報を格納し、単語リストセクションに単語リストを格納する、コンピュータによって実施される方法を含み、単語リストは、選択された単語に関連付けられた対応するポインタを識別する情報を有する。同様に、別の態様は、上述で提供されるレキシコンのデータ構造を使用して単語情報にアクセスすることである。

【 0 0 1 4 】

上述のレキシコン構造は、いくつかのレキシコンから情報を得ることが望ましい場合に特に有用であり、これは本発明のさらに別の態様である。一般に、各レキシコンのデータセクションに選択的にアクセスすることにより、特定の単語項目についての複数のレキシコンのデータを必要に応じて組み合わせる、無視する、または選択することができる。

【 0 0 1 5 】

【発明の実施の形態】

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。

【 0 0 1 6 】

図１に、通例はテキスト文字列の形で言語入力１２を受け取り、言語入力１２を処理して、通例は同じくテキスト文字列の形で言語出力１４を提供する言語またはテキスト処理システム１０を概略的に表す。数例を挙げると、言語処理システム１０は、例えばスペルチェッカ、文法チェッカ、あるいは自然言語プロセッサとして実施することができる。当業者には理解されるように、言語処理システム１０は、スタンドアロンアプリケーション、または別のシステムからのアクセスが可能な、あるいは別のシステムに含まれるモジュールやコンポーネントとすることができる。

【 0 0 1 7 】

一般に、言語処理システムはテキストアナライザ２０及びレキシコン２２を含む。テキス

10

20

30

40

50

トアナライザ 20 は、入力 12 を受け取り、レキシコン 22 にアクセスしてレキシコン 22 から情報を得、単語情報を処理して出力 14 を提供するコンポーネントまたはモジュールを図式的に表している。本発明の一態様は、その利用例によって必要とされる可能性のある必要な情報を効率的にテキストアナライザ 20 に提供するためのレキシコン 22 の改良されたデータ構造である。レキシコン 22 は、多くの言語処理システムと多くの形態のテキストアナライザに使用できる独立したコンポーネントであることを考慮して、テキストアナライザ 20 とレキシコン 22 の一般的な相互作用について説明するが、本発明の理解に必要でないため、各種形態のテキストアナライザに関する具体的な詳細については述べない。

【0018】

本発明のさらなる詳細な説明に入る前に、動作環境の概要を見ておくとうるである。図 2 は、本発明の実施が可能な適切なコンピューティングシステム環境 50 の一例である。コンピューティングシステム環境 50 は適切なコンピューティング環境の一例に過ぎず、本発明の使用または機能性の範囲について何らの制限を示唆するものではない。またコンピューティングシステム環境 50 は、その例示的動作環境に示す構成要素の任意の 1 つまたは組み合わせに関連する依存関係または必要性を有するものとも解釈すべきでない。

【0019】

本発明は、数多くの他の汎用または特殊目的のコンピューティングシステム環境または構成で動作することができる。本発明に使用するのに適している可能性があるよく知られるコンピューティングシステム、環境、及び/または構成の例には、これらに限定しないが、パーソナルコンピュータ、サーバコンピュータ、ハンドヘルドまたはラップトップデバイス、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 (personal computer)、ミニコンピュータ、メインフレームコンピュータ、上述のシステムまたはデバイスのいずれかを含む分散コンピューティング環境などがある。

【0020】

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には、特定タスクを行うか、あるい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、データ構造などが含まれる。本発明は、通信ネットワークを通じてリンクされたりモートの処理装置によってタスクを行う分散コンピューティング環境で実施することもできる。分散コンピューティング環境では、メモリ記憶装置を含むローカル及びリモート両方のコンピュータ記憶媒体にプログラムモジュールを置くことができる。これらのプログラム及びモジュールによって行われるタスクについて、以下で図面を用いて説明する。当業者は、以下の説明及び図を、任意形態のコンピュータ可読媒体に書き込むことのできるプロセッサ実行可能命令として実行することができる。

【0021】

図 2 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 60 の形態の汎用コンピューティングデバイスを含む。コンピュータ 60 の構成要素には、これらに限定しないが、プロセッサ 70、システムメモリ 80、及びシステムメモリを含む各種のシステムコンポーネントをプロセッサ 70 に結合するシステムバス 71 が含まれる。システムバス 71 は、各種のバスアーキテクチャの任意のものを使用したメモリバスまたはメモリコントローラ、ペリフェラルバス、及びローカルバスを含む数タイプのバス構造のいずれでもよい。例として、このようなアーキテクチャには、ISA (Industry Standard Architecture) バス、MCA (Micro Channel Architecture) バス、EISA (Enhanced ISA) バス、VESA (Video Electronics Standards Association) バス、及びメザンバスとも称される PCI (Peripheral Component Interconnects) バスが含まれるが、これらに限定しない。

【0022】

コンピュータ 60 は通例各種の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 60 がアクセスできる任意の利用可能媒体でよく、揮発性及び不揮発性の媒体、リムーバル及びノンリムーバル媒体が含まれる。これに限定しないが、例としてコンピュータ可読媒体は、コンピュータ記憶媒体及び通信媒体を含むことができる。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、またはその他のデータなどの情報を記憶するための任意の方法または技術で実現された不揮発性及び不揮発性の媒体、リムーバル及びノンリムーバルの媒体を含む。コンピュータ記憶媒体には、これらに限定しないが、R A M (random access memory)、R O M (read only memory)、E E P R O M (electrically erasable PROM)、フラッシュメモリまたは他のメモリ技術、C D (compact disc [disk]) - R O M、デジタル多用途ディスク (D V D) 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、あるいは、所望の情報の記憶に用いることができ、コンピュータ 60 によるアクセスが可能な任意の他の媒体が含まれる。

10

【0023】

通信媒体は、通例、搬送波または他の搬送機構などの、変調データ信号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統合し、任意の情報伝達媒体を含む。用語「変調データ信号」とは、信号中に情報を符号化するような方式でその特性の1つまたは複数を設定または変化させた信号を意味する。例として、通信媒体には、有線ネットワークまたは直接配線接続などの有線媒体と、音響、R F (radio frequencies)、赤外線、及び他の無線媒体などの無線媒体が含まれるが、これらに限定しない。上記の媒体のいずれの組み合わせもコンピュータ可読媒体の範囲に含めるべきである。

20

【0024】

システムメモリ 80 には、R O M 81 及び R A M 82 など、揮発性及び/または不揮発性メモリの形態のコンピュータ記憶媒体が含まれる。起動時などにコンピュータ 60 内の要素間の情報転送を助ける基本ルーチンを含む B I O S (Basic Input/Output System) 83 は、通例 R O M 81 に記憶される。R A M 82 は通例、プロセッサ 70 から即座にアクセス可能な、かつ/または現在プロセッサ 70 によって操作中のデータ及び/またはプログラムモジュールを含む。これらに限定しないが、例として、図 2 にはオペレーティングシステム 84、アプリケーションプログラム 85、他のプログラムモジュール 86、及びプログラムデータ 87 を示している。

30

【0025】

コンピュータ 60 は、他のリムーバル/ノンリムーバル、揮発性/不揮発性のコンピュータ記憶媒体も含むことができる。単なる例として、図 2 には、ノンリムーバル、不揮発性の磁気媒体の読み取りまたは書き込みを行うハードディスクドライブ 91、リムーバル、不揮発性の磁気ディスク 102 の読み取りまたは書き込みを行う磁気ディスクドライブ 101、及び C D - R O M や他の光媒体などのリムーバル、不揮発性の光ディスク 106 の読み取りまたは書き込みを行う光ディスクドライブ 105 を示す。例示的動作環境で使用できるこの他のリムーバル/ノンリムーバル、揮発性/不揮発性のコンピュータ記憶媒体には、これらに限定しないが、磁気テープカセット、フラッシュメモリカード、デジタル多用途ディスク、デジタルビデオテープ、ソリッドステート R A M、ソリッドステート R O M などが含まれる。ハードディスクドライブ 91 は通例、インタフェース 90 などのノンリムーバルのメモリインタフェースを通じてシステムバス 71 に接続され、磁気ディスクドライブ 101 及び光ディスクドライブ 105 は通例、インタフェース 100 などリムーバルなメモリインタフェースによってシステムバス 71 に接続される。

40

【0026】

上記で説明し、図 2 に示したドライブとそれに関連付けられた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、及びコンピュータ 60 のその他のデータの記憶を提供する。例えば図 2 では、ハードディスクドライブ 91 にオペレーティングシステム 94、アプリケーションプログラム 95、他のプログラムモジュール 96、及びプログラムデータ 97 を記憶している。これらのコンポーネントは、オペレー

50

ティングシステム 8 4、アプリケーションプログラム 8 5、他のプログラムモジュール 8 6、及びプログラムデータ 8 7と同じものでも、異なるものでもよいことに留意されたい。ここではオペレーティングシステム 8 4、アプリケーションプログラム 8 5、他のプログラムモジュール 8 6、及びプログラムデータ 8 7には、それらが少なくとも異なるコピーであることを表すために異なる参照符号をつけている。

【 0 0 2 7 】

ユーザは、キーボード 1 1 2、マイクロフォン 1 1 3、手書きタブレット 1 1 4、及びマウス、トラックボール、タッチパッドなどのポインティングデバイス 1 1 1 などの入力装置を通じてコンピュータ 6 0 にコマンドと情報を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。これら及び他の入力装置は、システムバスに結合されたユーザ入力インタフェース 1 1 0 を通じてプロセッサ 7 0 に接続することが多いが、パラレルポート、ゲームポート、あるいはユニバーサルシリアルバス（USB）など他のインタフェース及びバス構造によって接続することも可能である。モニタ 1 4 1 または他タイプの表示装置も、ビデオインタフェース 1 4 0 などのインタフェースを介してシステムバス 7 1 に結合される。コンピュータは、モニタに加えて、スピーカ 1 4 7 やプリンタ 1 4 6 など他の周辺出力装置も含むことができ、それらは出力周辺インタフェース 1 4 5 を通じて接続することができる。

【 0 0 2 8 】

コンピュータ 6 0 は、リモートコンピュータ 1 3 0 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する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 1 3 0 はパーソナルコンピュータ、ハンドヘルドデバイス、サーバ、ルータ、ネットワーク PC、ピアデバイス、あるいはその他の一般的なネットワークノードでよく、通例はコンピュータ 6 0 との関連で上記で挙げた要素の多くまたはすべてを含む。図 2 に示す論理接続には、構内ネットワーク（LAN）1 2 1 と広域ネットワーク（WAN）1 2 3 が含まれるが、他のネットワークを含むことも可能である。このようなネットワーキング環境は、オフィス、企業内のコンピュータネットワーク、イントラネット、及びインターネットに一般的に見られる。

【 0 0 2 9 】

LAN ネットワーキング環境で使用する場合、コンピュータ 6 0 はネットワークインタフェースまたはアダプタ 1 2 0 を通じて LAN 1 2 1 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 6 0 は通例、インターネットなどの WAN 1 2 3 を通じて通信を確立するためのモデム 1 2 2 またはその他の手段を含む。モデム 1 2 2 は内蔵型でも外付け型でもよく、ユーザ入力インタフェース 1 1 0 または他の適切な機構を介してシステムバス 7 1 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 6 0 との関連で図示したプログラムモジュール、またはその一部はリモートのメモリ記憶装置に格納することができる。これに限定しないが、例として図 2 ではリモートアプリケーションプログラム 1 3 5 がリモートコンピュータ 1 3 0 に常駐している。図のネットワーク接続は例示的なものであり、コンピュータ間に通信リンクを確立する他の手段を使用できることは理解されよう。

【 0 0 3 0 】

テキストアナライザ 2 0 は、コンピュータ 6 0、またはリモートコンピュータ 1 3 0 などコンピュータ 6 0 と通信する任意のコンピュータに常駐できることを理解されたい。同様に、レキシコン 2 2 は、コンピュータ 6 0 の上述の記憶装置の任意のものに常駐するか、または適切な通信リンクを通じてアクセス可能にすることができる。

【 0 0 3 1 】

図 3 は、レキシコン 2 2 の図式表現である。図の例示的实施形態では、レキシコン 2 2 は、ヘッダセクション 1 6 0、単語リストセクション 1 6 2、索引テーブルセクション 1 6 4、索引セクション 1 6 6、2 つ以上のレキシコンデータセクション 1 6 8（ここでは例として 1 6 個のセクション、1 6 8 a、1 6 8 b、1 6 8 c、1 6 8 d、1 6 8 e、1 6

10

20

30

40

50

8 f、1 6 8 g、1 6 8 h、1 6 8 i、1 6 8 j、1 6 8 k、1 6 8 l、1 6 8 m、1 6 8 n、1 6 8 o、1 6 8 p)、及び文字列ヒープセクション 1 7 0を含む。

【 0 0 3 2 】

ヘッダセクション 1 6 0 は一般に、レキシコン 2 2 の構造についての情報を格納する。ヘッダセクション 1 6 0 は、例えば、レキシコンの名前やバージョンについての情報を含むことができる。ヘッダセクション 1 6 0 はまた、メモリオフセット及び各セクション 1 6 2、1 6 4、1 6 6、1 6 8 a ~ 1 6 8 p、及び 1 7 0 のサイズについての情報も含むことができる。セクション 1 6 2 は、レキシコン 2 2 の単語リストを含む。セクション 1 6 2 に単語リストを実施するには、任意の適切なフォーマットを用いることができる。特に有用なフォーマットの 1 つは、よく知られるデータ構造技術である「トライ (t r i e)」構造で単語リストを格納するものである。このフォーマットの利点としては、特定の接頭辞で始まる可能性のある単語がいくつあるかを容易に判定できることが挙げられ、これは例えば手書き認識や、ユーザが特定の文字を書いた可能性を確かめる必要がある際に有用である。このフォーマットではまた、トラバース (traversal) のパスを前方向と後ろ方向の両方で知ることができる。上記のように、セクション 1 6 2 には他の形態の単語リストイングを使用することができる。例えば単純なテーブルやリストを使用することができる。さらに別の実施形態では「差分」技術を使用して単語リストを格納することができ、この場合は連続した単語の記号または文字の違いを格納する。

10

【 0 0 3 3 】

セクション 1 6 4 を説明する前に、セクション 1 6 6 と、複数のセクション 1 6 8 とのその関係を初めに説明しておくとう有用であろう。「従来の技術」の項で述べたように、現在のレキシコンでは、必要とするのが情報の一部だけであっても、特定の単語項目に関連付けられたすべての情報を読み出すことが必要とされる。セクション 1 6 8 a ~ 1 6 8 p は、レキシコン中の各単語項目のデータを所望の方式で編成することを可能にし、関連するレキシコン情報を概ねともにグループ化することができる。例えば、セクション 1 6 8 a ~ 1 6 8 p の 1 つを使用してスペルチェックに関連する情報を格納し、別のセクションに標準的な言語分類に関連する情報を格納することができる。索引セクション 1 6 6 は一般に、単語リストセクション 1 6 2 の単語項目に応じて、セクション 1 6 8 a ~ 1 6 8 p に格納されたデータへのポインタ (例えばセットとしてグループ化した) を提供する。すなわち、単語リストセクション 1 6 2 (例えばトライ構造) は、索引セクション 1 6 6 へのアクセスポイント (オフセット) を直接または間接的に決定する。一般に、単語情報を得るこの方法は、所与の単語に応じて単語リストセクションにアクセスして、索引セクションのポインタ識別を確定することを含む。このポインタ識別を使用して、索引セクションで単語のポインタを得る。次いでこのポインタを使用して複数のデータセクションのうちのどのデータセクションに所与の単語についての情報があるか、そしてその情報がそのデータセクションのどこに位置するかを確定する。したがって、セクション 1 6 2 にある特定の単語項目に対して、索引セクション 1 6 6 を通じて、セクション 1 6 8 a ~ 1 6 8 p に格納されたその単語の対応するレキシコンデータに選択的にアクセスすることができ、それにより所与の単語のすべての単語情報を処理するか、または少なくとも読み出す必要がない。

20

30

40

【 0 0 3 4 】

特に有用な一実施形態では、セクション 1 6 2 にある各単語項目についてのセクション 1 6 8 a ~ 1 6 8 p の索引セクション 1 6 6 中のポインタまたはポインタのセットを、その単語項目が名詞か、動詞か、形容詞かなど、その品詞 (「 P O S (part of speech) 」) によって分類する。したがって、ある単語項目の P O S についてのデータは、セクション 1 6 8 a ~ 1 6 8 p の P O S 情報への一連のポインタとなる。したがって、2 つの P O S を持つ単語項目には、セクション 1 6 6 に 2 つの別個のポインタセットがある。セットの 1 つは、第 1 の P O S (例えばその項目の名詞形) についての情報の位置を示し、第 2 のポインタセットは、もう一方の P O S (例えばその項目の動詞形) についての情報の位置を示す。このように索引セクション 1 6 6 は、単語項目の P O S に基づいた、あるレベル

50

のレキシコンデータ分類を提供する。ここで、レキシコン 22 が対象とする言語に応じて、POS 以外の他の分類形態を使用できることを理解されたい。例えば日本語や中国語には、品詞の代わりに屈折または声調による分類を使用することができる。ここでは索引セクション 166 が POS 分類を提供するものと例示しているが、この機能は制限的あるいは必須とみなすべきでない。

【0035】

また、中国語や日本語のような言語で使用する際には、ここで用いる意味の単語「単語 (word)」の使用には、記号、表意文字、語標 (logogram) などにも含まれることにも留意されたい。したがって、本発明の態様を使用してこれらの言語のレキシコンも構築することができ、そのレキシコンは特に断らない限りは特許請求の範囲に包含されるものとする。

10

【0036】

例示的实施形態では、各ポインタは、そのポインタがセクション 168a ~ 168p のうちどれをポイントするのかに関する情報、POS の種類に関連付けられた情報、識別されたセクション 168a ~ 168p 中で関連するデータを見つけられるオフセット値を含む。セクション 162 の所与の単語項目に関連付けられたポインタは固定することができるが、例示的实施形態では、各単語項目のポインタ数は単語項目ごとに異なってよい。このようにして、本質的な制限が常に伴うことなく、索引セクション 166 をより小型かつ柔軟にすることができる。

【0037】

20

ある単語項目のセクション 166 の例示的ポインタの概略表現を次に示す。

$X_1 : X_2 : X_3 : X_4$

ここで X_1 は単語項目のポインタセットの最後を示すフラグであり、 X_2 はセクション 168 の 1 つを識別する情報であり、 X_3 は POS または他の分類を識別する情報であり、 X_4 は、 X_2 によって識別されるレキシコンデータのオフセットを示す値である。このフォーマットを使用して、所与の単語についてのすべての情報のポインタを連続して格納することができ、単語リストセクションに応じて直接あるいは間接的に第 1 のポインタが識別され、最後のポインタのフラグ X_1 をセットして所与の単語のポインタリストの最後を示す。一実施形態では、索引セクション 166 は大きな D W O R D 配列 (迅速なアクセスのために 4 バイトの量、4 バイトワード配置) である。この実施形態では、1 バイトの内訳は、ポインタセットの最後を示す X_1 の 1 ビットフラグ、セクション 168a ~ 168p を示す 4 ビットの X_2 、及び POS の種類を示す 3 ビットの X_3 である。そして X_4 に 3 バイトを使用して、データがセクション 168a ~ 168p のどこに格納されているかを示す 24 ビットのオフセット値を提供する。このフォーマットは一例に過ぎず、他のフォーマットも使用できることを理解されたい。同様に、この例は必須あるいは制限的と解釈すべきでない。一般に、索引セクション 166 のポインタのフォーマットは、複数のセクション 168 中のデータの位置、及び必要な場合は単語情報の 1 つまたは複数の分類を示すために選択される。

30

【0038】

ここで、セクション 166 のポインタ項目のオフセット部分に入るのに十分な小ささのデータは、別個のセクション 168a ~ 168p ではなく、直接索引セクション 166 に符号化できることにも留意されたい。この種のデータの例にはスプリング情報や単語項目の確率及び頻度データが含まれるが、これらはいずれも多くの場合はデータオフセット値に割り当てられたビットに容易に格納することができる。

40

【0039】

上記のように、索引セクション 166 へのエントリは、セクション 162 の単語項目に応じて決まる。セクション 162 と 166 間の移行には各種の技術を使用することができる。第 1 の実施形態では、セクション 162 の各単語項目は必要とされるセクション 166 へのオフセットを含むことができる。ただしセクション 162 がトライ構造を備える場合は、トライの葉ノード構造の修正が必要となる場合がある。あるいは、トライ構造中のノ

50

ードのオフセットを索引セクション166へのオフセットとして使用することができる。例示的实施形態では、これは、ある単語項目についてのPOS索引のセットに40バイト(10個のPOSポインタ)を割り当てることを意味する。さらなる実施形態では、オフセット値をセクション166の単語項目の最後に付することができる。

【0040】

さらに別の実施形態では、索引テーブル164をレキシコン22の構造中に含める。索引テーブル164は単語項目とセクション166の索引とのマッピングを可能にし、これは、セクション166中のポインタの数が単語項目ごとに異なる可能性がある場合に特に有用である。ただし、関連付けられた単語項目ごとに、索引セクション166で固定サイズ数のポインタを使用することが可能である。この構造のセクション166を使用すると索引テーブルセクション164が不要になる。この代替実施形態では、セクション166の索引ポインタの固定数を超えるポインタを有する単語項目が許された場合は、オーバーフローテーブルを使用することができる。

10

【0041】

ここで、セクション162から索引セクション166へのオフセット、より具体的にはセクション168a~168p中のデータをポイントするセクション166のポインタは、レキシコン22からデータを検索する際の効率と速度を提供するように編成することができることに留意されたい。例えば、頻繁に使用される単語の他の情報の隣にセクション168a~168pのレキシコン情報を配置するようにオフセットポインタを編成することができ、あるいは、必要な場合は、セクション168a~168p中の関連付けられた情報同士をより近くに編成して、ハードディスク、フロッピー(登録商標)などのコンピュータ記憶装置に格納された際に情報検索時間を短縮することができる。

20

【0042】

セクション168a~168pのデータはその中に存在することができ、即ち、必要な場合は、同じセクション168a~168pに含まれる参照データへの、他のセクション168a~168pに含まれる参照データへの、セクション162のものと単語項目へのポインタを提供でき、及び例証的实施形態ではまた文字列ヒープ170へのポインタを提供することができる。文字列ヒープ170は、そのデータをセクション168a~168p中に複数の出現として格納する必要がある選択された文字列に単一の記憶位置を提供するために使用される。文字列ヒープ170は単一のセクションであるか、あるいはセクション168a~168pと同様のサブセクションを含むことができる。セクション168中の情報の他の形態には、決定木中のブルフラグ、値、単語リストなどがある。

30

【0043】

複数のセクション168a~168pを使用した単語項目データの編成により、それを実装するコンピュータのメモリを多量消費することなく、特定の用途の必要性を満たすようにレキシコン22を容易に適合することが可能になる。例えば、レキシコン22はRAMなどの高速アクセスメモリに読み込むことができるが、レキシコン中の特定タイプのデータが必要でない場合は、複数のセクション168a~168pのうちそのセクションを省略することができる。索引セクション166のポインタは、存在するセクション168a~168pだけを反映するように変更することができるが、さらなる実施形態では、セクション168a~168pが存在すれば情報が得られ、一方セクションが存在しなければ情報が求められないので、変更は不必要である。エラーが生じないように、レキシコンに存在するセクションは例えばヘッダ160に記録することができる。

40

【0044】

ここに記載するレキシコン構造の際立った利点は、入力12(図1)が単語を含む場合にレキシコンのユーザまたは作成者が後の検索のためにその単語についてのどのようなタイプの情報でも入れることができる点である。さらに、ユーザによって定義された情報はレキシコンに含まれる他の情報と混在させる必要がなく、複数のセクション168a~168pの専用のセクションに格納することができる。

【0045】

50

以下に挙げるのは、セクション 1 6 8 a ~ 1 6 8 p に適したセクションに編成されたレキシコンデータのいくつかの例である。これらは単なる例に過ぎず、レキシコン 2 2 のデータは、利便性または理解のために任意の所望の方式で編成できることに留意されたい。ここで述べるセクションは特に有用であることが判明しているが、必須あるいは制限的なものとは見なすべきではない。

【 0 0 4 6 】

形態データセクション - このような情報は、発音ならびに様々な単語の時制についてのその単語の他の形を含むことができる。

【 0 0 4 7 】

標準作成者データセクション - この情報には、その単語項目が単数形、複数形であるか否か、あるいはその単語が有生か無生かを示すデータを含むことができる。このセクションの単語項目に関連する情報は、一般にはその単語項目についてのよく知られた情報であり、素人でも作成することができる。このように、この情報はユーザの要件に合わせて用意に変更または修正することができる。

10

【 0 0 4 8 】

標準言語データセクション - この情報には単語項目の言語学的情報が含まれる。このような情報は一般の素人にはあまり知られないが、言語学者はこの情報を容易に理解し、必要に応じて修正することができる。

【 0 0 4 9 】

構文解析データセクション - この情報には、自然言語の構文解析に役立つ情報が含まれる。

20

【 0 0 5 0 】

領域 / 主題データセクション - この情報は領域または主題のコードに関連する。例えば、この情報により、対応する単語が物理、数学、地理、食物などに関連するものであることを示すことができる。

【 0 0 5 1 】

スペリングデータセクション - この情報はスペルチェック、例えば方言マーキング、制限マーキングなどに関連する。制限マーキングは、卑語、頭字語、古い語など許容されるがスペルチェックの際には提案されない単語を示す。

【 0 0 5 2 】

複数語表現データセクション - この情報は、イディオム、固有名、本や映画の題名、オフィスの名称、地名など、複数の単語を単独に識別する必要がある際に有用である。通例、各単語項目について格納されるデータは、複数単語表現でその単語の前にくる、かつ / または後に来る単語である。

30

【 0 0 5 3 】

例えば、複数のセクション 1 6 8 a ~ 1 6 8 p の 1 つが、そのレキシコン項目の作成者のみによって見つけられる、任意の階層的な名前値の対を含むことができる。例えば、作成者が複数語表現（上記）についての固有表現（N E ; Name Entity）情報を追加したい場合には、名前文字列に基づく値の対のセットをそのセクションに追加することができ、これを X M L 形式で表すと次のようになる。

40

```

<named-entity>
  <app-ne-id>movieFinder::the_longest_day</app-ne-id>
<semantic-type>movieFinder::movieTitle</semanticType>
  <genre>Drama</genre>
  <URL>http://www.movieFinder.com/fetch-movie-info/
    the_longest_day</URL>
<movie-info>
  <date>Jan.30, 1969</date>
  <running-time>137 min.</running-time>
  <studio>20th Century Fox</studio>
</movie-info>
<non-rated/>

```

10

```
</named-entity>
```

このように、このセクションは、任意の入れ子構造になった、文字列に基づく単純な値の名前の対を表すことができる。この形式はXMLのタグ属性をサポートせず、作成者がこのセクションの別個の下位要素としてその属性を符号化する。上の例では、映画の題名についてのデータは、必要な場合には格納することのできる利用例固有のデータの混合を含んでいる。

20

【0054】

このレキシコン22の構造は、レキシコンの各セクションはそのすぐ後のセクションと連続する必要がないという事実を利用することにより書き込み可能なレキシコンに対応する。すなわち、セクションは将来行われる拡張のために余分の未使用スペースを取っておくことができる。レキシコンへの更新操作は、該当する位置に新しい値を書き込むことによって行う。レキシコン22がDDL（ダイナミックリンクライブラリ）ベースのレキシコン、あるいは予備スペースを備えない事前にコンパイルされた（静的な）ファイルベースのレキシコンとして実施される場合は、単純なフリーリスト実装により、先頭一致（first-fit）アルゴリズムに基づいて空いている項目スペースを見つける。

30

【0055】

一般に、レキシコン22に単語情報を格納する方法は、複数のデータセクション168に単語情報を格納し、各データセクションは単語リスト中の単語について実質的に異なる選択された情報を格納することと、複数のデータセクション168と分離した索引セクション166にポインタ情報を格納し、各ポインタは複数のデータセクション168中の選択されたデータをポイントすることと、複数のデータセクション168及び索引セクション166と分離した単語リストセクション162に単語リストを格納し、単語リストは選択された単語に関連付けられた対応するポインタを識別する情報を有することとを含む。必要な場合は、識別値を索引テーブルセクション164に格納することができ、この場合には各識別値が単語リストセクション162の単語に対応し、索引セクション166のポインタと関連付けられる。同様に、ポインタに分類の表示を含めて単語情報を分類することもできる。

40

【0056】

このレキシコン22の構造は、いくつかのレキシコンから情報を得ることが望ましい場合に特に有用である。一般には、特定の単語項目についての複数のレキシコンの情報を必要に応じて組み合わせる、無視する、あるいは選択することができる。いくつかのレキシコンからのレキシコン情報を組み合わせる例は、核または基礎となるレキシコンに単語項目

50

についての第1の量の情報が含まれ、第2のレキシコンに特定領域についての単語項目についての第2の量の情報が含まれ、第3のレキシコンにユーザが決定した単語項目についての第3の量の情報が含まれる実装に見られる。

【0057】

図4に、特定の単語項目についての情報を複数のレキシコンから得る方式を図式的に示している。図4で、レキシコン（データセクションだけによって表しているが、本来は図3に示すセクションの一部またはすべてを備える）は行に編成され、これを180、181、182、及び183で示す。図4では個々のデータセクション（168に対応する）を縦に表しており、この例証的实施形態では、4つのレキシコン180～183を通じて最大6個のデータセクション190、191、192、193、194、及び195にアクセスすることができる。各レキシコン180～183がデータセクション190～195すべてを含むことは必須でなく、多くの実際の事例では、すべてのレキシコンのすべてのデータセクション間にそのような対応関係は存在しないことに留意されたい。

【0058】

図4では符号 X_y を使用してレキシコンセクションのデータを示しており、 X はデータセクション190～195の1つを現し、 Y はレキシコン180～183を表す。例えばレキシコン180は、データセクション190₁₈₀、193₁₈₀、及び195₁₈₀を備える。

【0059】

レキシコン180～183のデータは同じタイプの内容を有するセクション190～195に編成されているので、レキシコン180～183にまたがって情報を容易に組み合わせる、あるいは選択することができる。所与の単語項目について、第1のレキシコン180の情報を調べ、次いで必要に応じて他のレキシコン181～183の同じデータセクションに進むことにより情報を得ることができる。一実施形態では、検索するデータは、実行時に定義される所望のセクションタイプのセットによって制御する。1つの変数で、レキシコンのある項目のデータを読み出すか読み出さないかを決定する。第2の変数で、調べた他のレキシコンの対応するセクションから先に読み出された項目のデータと組み合わせる、あるいはそのデータに上書きするかどうかを決定する。概略的には、レキシコンは「スタック」されており、スタック中の最上位のレキシコンのデータセクション190～195を読み取り、次いでスタックを順次下に進み、読み出すか読み出さないか、選択するか、無視するか、上書きするか、あるいは組み合わせるかについての規則に従うことによって情報を得ると考えることができる。図3に示すレキシコン構造は、実装者が、所与のセクションタイプのデータを他のレキシコンの同じセクション中のデータと組み合わせる、あるいは上書きする方式を選択することを可能にする。

【0060】

図4で、レキシコン180～183から得た情報186は、セクション190₁₈₀、191₁₈₁、192₁₈₂、193₁₈₀、194₁₈₃、及び195₁₈₀₊₁₈₁₊₁₈₂に対応する情報を含む。この例では、セクション190、191、192、193、及び194のデータは、単にデータセクションの1つに停止インディケータが見つかるまで、レキシコン180～183をセクションごとに順に調べることによって得る。例えば、レキシコン180及び183はともにセクション190に情報を有するが、停止インディケータがセクション190₁₈₀で見つかるのでレキシコン180の情報だけが取り出される。実行時に、これによりセクション190₁₈₃の情報が無視される。これに対してセクション195₁₈₂を調べるまでに停止インディケータが見つからないので、セクション195₁₈₀、195₁₈₁、及び195₁₈₂の情報は組み合わせて情報195₁₈₀₊₁₈₁₊₁₈₂を形成する。必要な場合は、テキストアナライザ20、またはテキストアナライザ20の要求に基づいてレキシコン22にアクセスするインタフェースモジュール（図示せず）によって実施される規則に基づいて、すべてのレキシコンのセクションにわたる情報を組み合わせる、無視する、あるいはその他の形で選択することができる。そのような規則は、例えば、他のレキシコンの対応するセクションに情報がある

10

20

30

40

50

かどうかに関係なく、特定のレキシコンの特定のセクションを常に使用することを指定することができる。これを図4に表しており、ここではレキシコン182のセクション192に情報があり、この情報はスタック中で上方にあるので少なくとも最初は調べられるが、レキシコン183のセクション192の情報が得られる。ただしデータの選択は、例えば上述の要領で停止ポインタを使用して単語項目ごとにしてもよい。

【0061】

要約すると、以前には得られなかった柔軟性と効率性を提供する改良したレキシコン構造について述べた。索引セクション及び複数のデータセクションにより、テキスト処理システム及び/または利用可能なコンピュータリソースの必要性に合わせてレキシコンを適合することが可能になる。この改良されたデータ構造により、複数のレキシコンのデータに選択的にアクセスし、かつ/または必要に応じて組み合わせることも可能になる。

10

【0062】

本発明について好適実施形態を参照して説明したが、当業者は、本発明の趣旨及び範囲から逸脱せずに形態及び詳細に変更を加えられることを認識されよう。

【0063】

【発明の効果】

以上説明したように本発明によれば、レキシコンを格納するための高い柔軟性と効率を有する改良されたデータ構造を提供しそのデータ構造を使用可能とすることができる。

【図面の簡単な説明】

【図1】本発明の実施形態の言語またはテキスト処理システムのブロック図である。

20

【図2】本発明の実施形態の例示的環境のブロック図である。

【図3】本発明の実施形態のレキシコンの図式表現の図である。

【図4】本発明の実施形態の複数のレキシコンにわたって情報を検索する、または情報にアクセスする図式表現の図である。

【符号の説明】

- 10 言語処理システム
- 12 入力
- 14 出力
- 20 テキストアナライザ
- 22 レキシコン
- 50 コンピューティングシステム環境
- 60、130 コンピュータ
- 70 プロセッサ
- 71 システムバス
- 80 システムメモリ
- 81 ROM
- 82 RAM
- 83 BIOS
- 84、94 オペレーティングシステム
- 85、95、135 アプリケーションプログラム
- 86、96 プログラムモジュール
- 87、97 プログラムデータ
- 90、100 インタフェース
- 91 ハードディスクドライブ
- 101 磁気ディスクドライブ
- 102 磁気ディスク
- 105、106 光ディスクドライブ
- 110 ユーザ入力インタフェース
- 111 ポインティングデバイス
- 112 キーボード

30

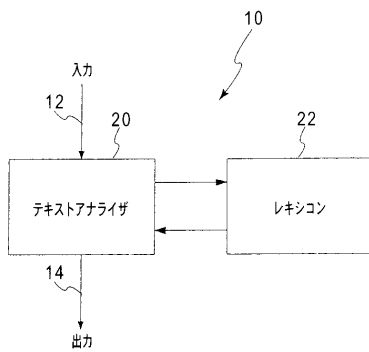
40

50

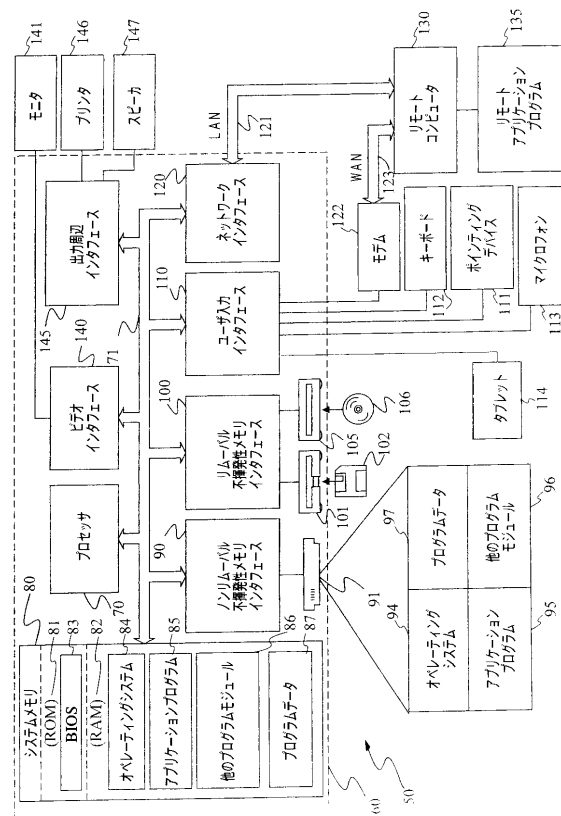
- 1 1 3 マイクロフォン
- 1 1 4 手書きタブレット
- 1 2 0 ネットワークインタフェース
- 1 2 1 LAN
- 1 2 2 モデム
- 1 2 3 WAN
- 1 4 0 ビデオインタフェース
- 1 4 1 モニタ
- 1 4 6 プリンタ
- 1 4 7 スピーカ
- 1 6 0、1 6 2、1 6 4、1 6 6、1 6 8、1 7 0
- 1 6 8 a ~ 1 6 8 p セクション
- 1 8 0、1 8 1、1 8 2、1 8 3 レキシコン
- 1 9 0、1 9 1、1 9 2、1 9 3、1 9 4、1 9 5 データセクション

10

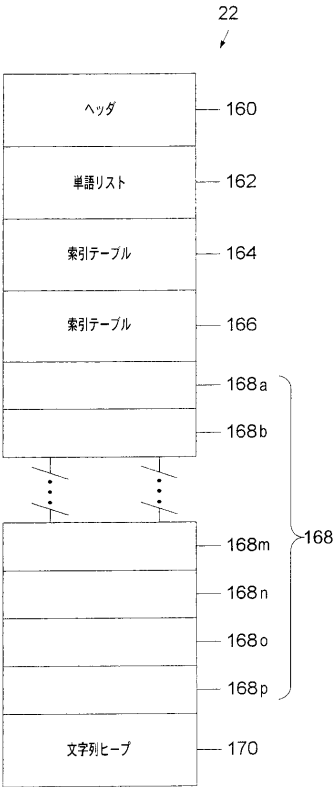
【図 1】



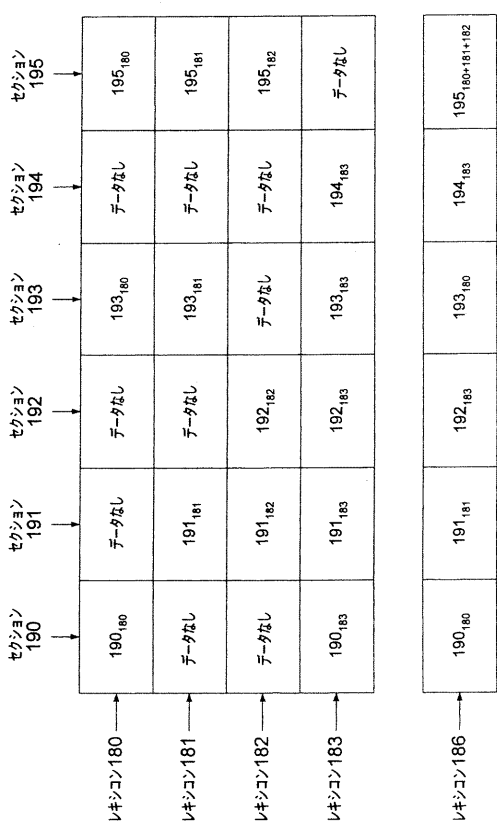
【図 2】



【図 3】



【図 4】



フロントページの続き

- (72)発明者 ジェームズ ピー・フィニガン
アメリカ合衆国 98007 ワシントン州 ベルビュー サウスイースト 6 14351 ア
パートメント 0201
- (72)発明者 カーティス イー・ハッテンハウアー
アメリカ合衆国 98007 ワシントン州 ベルビュー ザ レイクス ノースイースト 42
プレイス 14442 アパートメント 709
- (72)発明者 ダグラス ダブリュ・ポッター
アメリカ合衆国 98133 ワシントン州 シアトル ノース 128 ストリート 2155
- (72)発明者 ケビン アール・パウエル
アメリカ合衆国 98034 ワシントン州 カークランド ノースイースト 137 プレイス
13104

合議体

審判長 長島 孝志

審判官 久保 正典

審判官 飯田 清司

- (56)参考文献 特開平7-152756(JP,A)
特開平5-204962(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F17/28