

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4285272号
(P4285272)

(45) 発行日 平成21年6月24日(2009.6.24)

(24) 登録日 平成21年4月3日(2009.4.3)

(51) Int. Cl. F I
G06F 17/22 (2006.01) G O 6 F 17/22 5 2 2 L
G06F 17/27 (2006.01) G O 6 F 17/27 E

請求項の数 3 (全 17 頁)

<p>(21) 出願番号 特願2004-53456 (P2004-53456) (22) 出願日 平成16年2月27日 (2004.2.27) (65) 公開番号 特開2005-242809 (P2005-242809A) (43) 公開日 平成17年9月8日 (2005.9.8) 審査請求日 平成18年10月31日 (2006.10.31)</p> <p>前置審査</p>	<p>(73) 特許権者 000002945 オムロン株式会社 京都市下京区塩小路通堀川東入南不動堂町 801番地</p> <p>(74) 代理人 100078916 弁理士 鈴木 由充</p> <p>(74) 代理人 100142114 弁理士 小石川 由紀乃</p> <p>(72) 発明者 伴 哲也 京都府京都市下京区西洞院木津屋橋通東入 ル オムロンソフトウェア株式会社内</p> <p>(72) 発明者 山田 堂弘 京都府京都市下京区西洞院木津屋橋通東入 ル オムロンソフトウェア株式会社内</p> <p style="text-align: right;">最終頁に続く</p>
--	--

(54) 【発明の名称】 形態素解析方法、この方法を用いたプログラムおよび情報処理装置

(57) 【特許請求の範囲】

【請求項1】

所定数の変換後文字列が登録された変換用標準辞書および新規の変換後文字列を登録するための学習用辞書、ならびに所定数の形態素が登録された形態素解析用辞書が格納されたメモリを具備し、操作部よりかな文字列の入力を受け付けた後に辞書検索に基づく変換処理を指示する操作を受け付けたことに応じて、入力されたかな文字列により変換用標準辞書および学習用辞書を検索して、変換後文字列の候補を抽出する機能と、前記操作部よりかな文字列の入力を受け付けた後に、カタカナ、ひらがな、アルファベット、数字のいずれか一の文字種を選択して変換を指示する後変換操作を受け付けたことに応じて、前記入力されたかな文字列を選択された文字種による文字列に変換する機能と、前記後変換操作に応じて変換され、前記変換用標準辞書および学習用辞書に登録されていない変換後文字列を前記学習用辞書に登録する機能とを具備する文字変換手段として動作するコンピュータにおいて、漢字およびかな混じりの文字列を含むテキストデータの入力を受け付けてそのテキストデータを形態素に分解する形態素解析処理を実行する方法であって、

前記形態素解析処理の対象となるテキストデータの入力を受け付けたとき、当該テキストデータから所定長さの文字列を抽出するステップAと、ステップAで抽出された文字列により前記形態素解析用辞書を検索して形態素の候補を抽出するステップBと、前記ステップAで抽出された文字列を前記文字変換手段に渡して当該文字列による学習用辞書の検索を実行させ、その検索により抽出された文字列を形態素の候補として前記文字変換手段より受け付けるステップCとを複数サイクル実行した後に、各サイクルのステップBおよ

びステップCにより得た候補の中の所定数を選択して出力するステップDを実行し、
前記ステップCでは、前記ステップAで抽出された文字列がひらがな文字列である場合
には、この文字列を前記文字変換手段に渡さずに、形態素の候補を得られなかったという
処理結果を設定し、

前記ステップDでは、前記ステップCにより得たカタカナ文字列、アルファベット文字
列、および数字による文字列を、ステップBにより得た候補より優先して選択する、
ことを特徴とする形態素解析方法。

【請求項2】

所定数の変換後文字列が登録された変換用標準辞書および新規の変換後文字列を登録す
るための学習用辞書が格納されたメモリを具備し、操作部よりかな文字列の入力を受け付
けた後に辞書検索に基づく変換処理を指示する操作を受け付けたことに応じて、入力され
たかな文字列により変換用標準辞書および学習用辞書を検索して、変換後文字列の候補を
抽出する機能と、前記操作部よりかな文字列の入力を受け付けた後に、カタカナ、ひらが
な、アルファベット、数字のいずれか一の文字種を選択して変換を指示する後変換操作を
受け付けたことに応じて、前記入力されたかな文字列を選択された文字種による文字列に
変換する機能と、前記後変換操作に応じて変換され、前記変換用標準辞書および学習用辞
書に登録されていない変換後文字列を前記学習用辞書に登録する機能とを具備する文字変
換手段として動作するコンピュータに、形態素解析処理の機能を付与するためのプログラ
ムであって、

漢字およびかな混じりの文字列を含むテキストデータの入力を受け付ける第1のステッ
プ；および入力されたテキストデータを形態素に分解して各形態素の組み合わせを出力す
る第2のステップを前記コンピュータに実行させるためのプログラムと、所定数の形態素
が登録された形態素解析用辞書を構成する電子データとを含み、

前記第2のステップでは、前記第1のステップで受け付けたテキストデータから所定長
さの文字列を抽出するステップAと、ステップAで抽出された文字列により前記形態素解
析用辞書を検索して形態素の候補を抽出するステップBと、前記ステップAで抽出された
文字列を前記文字変換手段に渡して当該文字列による学習用辞書の検索を実行させ、その
検索により抽出された文字列を形態素の候補として前記文字変換手段より受け付けるステ
ップCとを複数サイクル実行した後に、各サイクルのステップBおよびステップCにより
得た候補の中の所定数を選択して出力するステップDを実行し、

前記ステップCでは、前記ステップAで抽出された文字列がひらがな文字列である場合
には、この文字列を前記文字変換手段に渡さずに、形態素の候補を得られなかったという
処理結果を設定し、

前記ステップDでは、前記ステップCにより得たカタカナ文字列、アルファベット文字
列、および数字による文字列を、ステップBにより得た候補より優先して選択する、
ことを特徴とする形態素解析用のプログラム。

【請求項3】

操作部と、この操作部よりかな文字列の入力および変換指示操作を受け付けて、入力さ
れたかな文字列を他の形態の文字列に変換する文字変換手段と、漢字およびかな混じりの
文字列を含むテキストデータの入力を受け付けて、そのテキストデータを形態素に分解す
る形態素解析手段とを具備する装置であって、

前記文字変換手段は、

所定数の変換後文字列が登録された変換用標準辞書と、

新規の変換後文字列を登録するための学習用辞書と、

前記操作部よりかな文字列の入力を受け付けた後に辞書検索に基づく変換処理を指示す
る操作を受け付けたことに応じて、入力されたかな文字列により変換用標準辞書および学
習用辞書を検索して、変換後文字列の候補を抽出する候補検索手段と、

前記操作部よりかな文字列の入力を受け付けた後に、カタカナ、ひらがな、アルファベ
ット、数字のうちのいずれか一の文字種を選択して変換を指示する後変換操作を受け付
けたことに応じて、前記入力されたかな文字列を指示された文字種による文字列に変換する

10

20

30

40

50

後変換操作時処理手段と、

前記後変換操作時処理手段により変換され、前記変換用標準辞書および学習用辞書に登録されていない変換後文字列を前記学習用辞書に登録するとともに、前記形態素解析手段から検索対象の文字列を受け付けたとき、この文字列により学習用辞書を検索して、その検索により抽出した変換後文字列を形態素解析手段に渡す学習用辞書処理手段とを具備し

、
前記形態素解析手段は、

所定数の形態素が登録された形態素解析用辞書と、

処理対象のテキストデータのを受け付けたとき、当該テキストデータから所定長さの文字列を抽出するステップAと、ステップAで抽出された文字列により前記形態素解析用辞書を検索して形態素の候補を抽出するステップBと、前記ステップAで抽出された文字列を前記文字変換手段の学習用辞書処理手段に渡し、当該文字列による学習用辞書の検索を実行した学習用辞書処理手段から渡された文字列を形態素の候補として受け付けるステップCとを複数サイクル実行する検索手段と、

前記検索手段による各サイクルのステップBおよびステップCにより得た候補の中の所定数を選択して出力する出力手段とを、具備し、

前記検索手段は、前記ステップAで抽出された文字列がひらがな文字列である場合には、ステップCにおいて、前記抽出された文字列を前記学習用辞書処理手段に渡さず、形態素の候補を得られなかったという処理結果を設定し、

前記出力手段は、前記ステップCにより得たカタカナ文字列、アルファベット文字列、および数字による文字列を、ステップBにより得た候補より優先して選択する、情報処理装置。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、コンピュータに所定容量のテキストデータを入力して、そのテキストデータに対する形態素解析を行う技術に関する。

【背景技術】

【0002】

近年、インターネットの普及に伴い、ウェブページや電子メールなどに含まれるテキストデータに形態素解析を行って、抽出された単語を他の文書を作成する目的で使用することが提案されている。また、前記形態素解析の結果に基づき、テキストデータを音声データに変換して出力するようにしたソフトウェアなども開発されている。

【0003】

一般的な形態素解析処理では、種々の形態素が登録された辞書を用いて処理対象のテキストデータに含まれる文字列を検索し、複数の候補を抽出する。さらに、形態素の組み合わせにかかる規則に基づいて最適な候補の組み合わせを決定し、その決定による形態素の組み合わせを出力する。

このような形態素解析において、解析の精度を向上するには、辞書に登録されていない未知語を抽出できるようにする必要がある。この点につき、下記の特許文献1のような先行技術が存在する。

【0004】

【特許文献1】特開平6-12453号 公報

【0005】

上記の特許文献1では、未登録の単語やその前後の結合関係を種々のルールと照合することにより未知語を抽出した後、抽出された未知語をモニタに表示する。そして、ユーザーに登録すべき未知語を選択させた上で、読みや品詞情報などのを受け付け、これらに対応づけた新たな辞書データを作成してメモリに登録するようにしている。

【発明の開示】

【発明が解決しようとする課題】

10

20

30

40

50

【 0 0 0 6 】

形態素解析では、世間の情勢や流行により生まれた新たな名称や言い回し、ユーザーが仲間うちで使用する単語など（以下、これらを「新語」と総称する。）に速やかに対応できるようにする必要がある。特に、近年は、タレントの名前、ブランド名、略語などを、カタカナ文字列やアルファベット文字列で表す頻度が高くなっているから、形態素解析でもこれらの新語に速やかに対応できるように、簡単な方法で新語を登録できるようにするのが望ましい。

【 0 0 0 7 】

上述した特許文献 1 に記載の発明では、新語を含むテキストデータを形態素解析した場合、未知語として抽出された新語について、ユーザーが登録作業を行わなければならない。このため、ユーザーの負担が大きくなるという問題がある。また、ユーザーが登録すべき新語を見落とししたり、登録作業に負担を感じて処理を中止した場合には、その新語は登録されないままとなるから、次の形態素解析で同じ新語を含むテキストデータを処理した場合、その新語は、再び未知語として取り扱われることになる。

このように、特許文献 1 の発明は、未知語の登録のために特別な処理時間が必要である上、ユーザーが未知語の登録作業を行うことを前提にしており、登録処理を簡単に行うのは困難である。

【 0 0 0 8 】

一方、システム提供者が定期的に形態素解析用辞書に新語を追加して、これをデータ配信などの方法で各ユーザーに提供するようにすれば、ユーザーの負担を軽減でき、形態素解析の機能を高めることもできる。しかし、新規登録すべき新語が多数になるのに対し、実際に使用される単語はそのうちの一部になる可能性が高く、しかも、使用される新語はユーザーによって異なるものになる可能性がある。

このように、形態素解析用辞書を定期的に更新する方法では、効率が悪く、システム提供者の採算に見合わないという問題が生じる。

【 0 0 0 9 】

ところで、一般ユーザーは、電子メールなどの文書を作成する際に、自分が関心を持つ新語を入力する可能性が高い。また、他者からの電子メールやウェブページなどにアクセスする場合にも、自分が関心を持つ新語が含まれるデータを閲覧する可能性が高いと考えられる。

【 0 0 1 0 】

かな漢字変換処理用のソフトウェアには、一般に、内部辞書に登録されていない文字列への変換が行われたときに、その変換後文字列を自動的に登録する学習機能が設定されている。また、この種のソフトウェアでは、かな漢字以外の特定の文字種（アルファベットやカタカナなど）への変換を指示する機能を所定のキーに割り当て、そのキーを操作することによって、入力文字列を特定の文字種による文字列に変換するようにしている。この特定の文字種の文字列への変換を指示するための操作は、「後変換操作」と呼ばれている。

【 0 0 1 1 】

この明細書では、特定の文字種への変換を指示する操作からその操作に応じて変換された文字列を確定する操作までを含めて、「後変換操作」という。

前記したカタカナやアルファベットにより表記される新語は、この後変換操作により入力することができる。また、この後変換操作により確定された文字列は、前記した学習機能により登録することができるので、次回、同じ新語を入力する際には、変換前のかな文字列を入力して通常の変換操作を行うことによって、前回登録された文字列を呼び出すことが可能となる。このように、文字入力処理では、ユーザーが特別の登録作業を行わなくとも、通常の方法で文字入力作業の過程で新語を辞書データに加えることが可能である。

【 0 0 1 2 】

この発明は、上記の点に着目し、文字変換処理で蓄積された辞書データを形態素解析処理でも使用できるようにすることにより、ユーザーが特別な登録作業を行ったり、形態素

10

20

30

40

50

解析用の辞書を定期的に更新しなくとも、形態素解析用の機能を自然に向上できるようにすることを目的とする。

また、この発明は、形態素解析において、個々のユーザーの関心に応じた新語を抽出できるようにすることによって、情報処理における利便性を高めることを目的とする。

【課題を解決するための手段】

【0013】

この発明にかかる形態素解析方法は、所定数の変換後文字列が登録された変換用標準辞書および新規の変換後文字列を登録するための学習用辞書、ならびに所定数の形態素が登録された形態素解析用辞書が格納されたメモリを具備し、操作部よりかな文字列の入力を受け付けた後に辞書検索に基づく変換処理を指示する操作を受け付けたことに応じて、入力されたかな文字列により変換用標準辞書および学習用辞書を検索して、変換後文字列の候補を抽出する機能と、前記操作部よりかな文字列の入力を受け付けた後に、カタカナ、ひらがな、アルファベット、数字のいずれか一の文字種を選択して変換を指示する後変換操作を受け付けたことに応じて、前記入力されたかな文字列を選択された文字種による文字列に変換する機能と、後変換操作に応じて変換され、前記変換用標準辞書および学習用辞書に登録されていない変換後文字列を前記学習用辞書に登録する機能とを具備する文字変換手段として動作するコンピュータにおいて実行されるもので、処理対象のテキストデータ（漢字およびかな混じりの文字列を含むテキストデータである。以下においても同じ。）の入力を受け付けたとき、当該テキストデータから所定長さの文字列を抽出するステップAと、ステップAで抽出された文字列により前記形態素解析用辞書を検索して形態素の候補を抽出するステップBと、ステップAで抽出された文字列を前記文字変換手段に渡して当該文字列による学習用辞書の検索を実行させ、その検索により抽出された文字列を形態素の候補として前記文字変換手段より受け付けるステップCとを複数サイクル実行した後に、各サイクルのステップBおよびCにより得た候補の中の所定数を選択して出力するステップDを実行する。また、ステップCでは、ステップAで抽出された文字列がひらがな文字列である場合には、この文字列を文字変換手段に渡さずに、形態素の候補を得られなかったという処理結果を設定する。またステップDでは、ステップCにより得たカタカナ文字列、アルファベット文字列、および数字による文字列を、ステップBにより得た候補より優先して選択する。

【0014】

上記において、文字変換手段により実行される変換処理では、入力されたかな文字列（ひらがな文字列である場合が多い。）を、漢字文字列のほか、アルファベット、数字、カタカナ、ひらがななどの特定の文字種による文字列に変換することができる。なお、この明細書でいうところの漢字文字列には、漢字のみから成る文字列のほか、漢字とかな（主としてひらがな）とを組み合わせた文字列も含まれるものとする。

【0015】

変換用標準辞書には、変換後文字列に読み／品詞／使用頻度などを対応づけたものを、辞書データとして格納することができる。学習用辞書にも同様の構成の辞書データを格納することができるが、この学習用辞書は最初は空の状態であり。また、これら2種類の辞書に加え、ユーザーの設定操作によって作成された辞書データを格納する辞書（ユーザー辞書）を設けてもよい。

【0018】

形態素解析用辞書には、形態素に読み／品詞／他の形態素との係り受けの関係などを対応づけた辞書データを格納することができる。形態素解析処理では、入力されたテキストデータを所定位置で区切ることにより検索キーワードとなる文字列を抽出し、その検索キーワードにより前記形態素解析用辞書および学習辞書を検索して、候補を抽出することができる。

【0019】

上記の形態素解析方法によれば、通常のかな文字列の変換処理を実行する際に行われた後変換操作に応じて学習用辞書に蓄積された変換後文字列を形態素解析処理でも使用する

10

20

30

40

50

ことができるので、特別の登録処理を行ったり、形態素解析用の辞書を更新しなくとも、形態素の解析機能を向上することができる。また、後変換操作によりかな文字列が名称や略語などの新語に変換され、その新語が学習用辞書に登録されると、形態素解析においても、登録された新語を含むテキストデータが入力された場合には、前記学習用辞書を用いた検索により、その新語を抽出することが可能となる。

また、この方法では、学習用辞書から抽出されたカタカナ文字列、アルファベット文字列、および数字による文字列を形態素解析用辞書から抽出された候補より優先的に選択するので、形態素解析用辞書に登録されていないカタカナ文字列、アルファベット文字列、および数字による文字列を、形態素解析処理で容易に抽出することが可能になる。

【0020】

つぎに、形態素解析処理のために受け付けたテキストデータから抽出された文字列がひらがな文字列である場合には、ステップCにおいて、この文字列を文字変換手段に渡さずに形態素の候補を得られなかったという処理結果を設定するのは、ひらがな文字列は、『送りがな』などを特定する用途で入力されることが多く、重要な意味のある単語を表す可能性が低いためである。このようにすれば、後変換操作により学習用辞書に登録された文字列のうち、名称や略語を表す可能性が高いカタカナ文字列、アルファベット文字列、数字による文字列のみが形態素解析処理で抽出されるようにすることができ、検索にかかる時間を短縮することができる。

【0022】

つぎに、この発明にかかる形態素解析用のプログラムは、所定数の変換後文字列が登録された変換用標準辞書および新規の変換後文字列を登録するための学習用辞書が格納されたメモリを具備し、操作部よりかな文字列の入力を受け付けた後に辞書検索に基づく変換処理を指示する操作を受け付けたことに応じて、入力されたかな文字列により変換用標準辞書および学習用辞書を検索して、変換後文字列の候補を抽出する機能と、前記操作部よりかな文字列の入力を受け付けた後に、カタカナ、ひらがな、アルファベット、数字のいずれか一の文字種を選択して変換を指示する後変換操作を受け付けたことに応じて、前記入力されたかな文字列を選択された文字種による文字列に変換する機能と、後変換操作に応じて変換され、前記変換用標準辞書および学習用辞書に登録されていない変換後文字列を前記学習用辞書に登録する機能とを具備する文字変換手段として動作するコンピュータに導入される。

【0023】

上記のプログラムは、処理対象のテキストデータの入力を受け付ける第1のステップ；入力されたテキストデータを形態素に分解して各形態素の組み合わせを出力する第2のステップ；の各ステップをコンピュータに実行させるためのプログラムと、前記した形態素解析用辞書を構成する電子データとを含む。このうちの第2のステップでは、前記第1のステップで受け付けたテキストデータから所定長さの文字列を抽出するステップAと、ステップAで抽出された文字列により形態素解析辞書を検索して形態素の候補を抽出するステップBと、ステップAで抽出された文字列を文字変換手段に渡して当該文字列による学習用辞書の検索を実行させ、その検索により抽出された文字列を形態素の候補として前記文字変換手段より受け付けるステップCとを複数サイクル実行した後に、各サイクルのステップB及びステップCにより得た候補の中の所定数を選択して出力するステップDを実行する。

【0024】

上記の形態素解析用のプログラムでは、ステップCにおいて、ステップAで抽出された文字列がひらがな文字列である場合には、この文字列を文字変換手段に渡さずに、形態素の候補を得られなかったという処理結果を設定する。またステップDでは、ステップCにより得たカタカナ文字列、アルファベット文字列、および数字による文字列を、ステップBにより得た候補より優先して選択する。

【0026】

上記のプログラムによれば、後変換操作により入力された単語を抽出可能な形態素解析

10

20

30

40

50

処理を実行することができる。また、このプログラムを一度インストールすれば、形態素解析用辞書の更新処理を行わなくとも、後変換操作に応じて学習用辞書に蓄積された変換後文字列を形態素解析で使用することができる。すなわち、ユーザーは、形態素解析のために特別な登録作業を行う必要も、形態素解析用辞書の更新データをインストールする必要もなしに、通常の文字変換処理を行うだけで、形態素解析の能力を向上させることができる。また、システム開発者も、形態素解析用辞書を定期的に更新することなく、各ユーザーの要望に応えることができる。

【0027】

この形態素解析用のプログラムは、ウェブブラウザやメールリーダーなどのアプリケーション（以下、「上位アプリケーション」という。）の稼働時に、ユーザーの操作などに応じて起動させることができる。この場合のステップDでは、前記決定した形態素の組み合わせを、前記上位アプリケーションに出力することができる。また、これら上位アプリケーションとは別のアプリケーション（たとえば音声出力システム）に、決定した形態素の組み合わせを渡すこともできる。

【0028】

上記のプログラムは、パーソナルコンピュータのほか、携帯電話、PDAなどの携帯端末の制御部を構成するコンピュータに組み込むことができる。また、このプログラムは、CD-ROMなどの記憶媒体に格納する方法や、電気通信回線により伝送する方法によって、ユーザーに提供することができる。また、この形態素解析用のプログラムと前記文字変換処理用のプログラムとを組み合わせたものを、1つのパッケージソフト（上記2種類のプログラムが格納された1または複数の記憶媒体から成り、各プログラムをコンピュータに同時または選択的にインストールできるようにしたもの）として提供したり、電気通信回線を介して所定のコンピュータに提供することができる。

【0029】

つぎに、この発明にかかる情報処理装置は、操作部と、この操作部よりかな文字列の入力および変換指示操作を受け付けて、入力されたかな文字列を他の形態の文字列に変換する文字変換手段と、処理対象のテキストデータの入力を受け付けて、そのテキストデータを形態素に分解する形態素解析手段とを具備するものである。前記文字変換手段は、所定数の変換後文字列が登録された変換用標準辞書と、新規の変換後文字列を登録するための学習用辞書と、前記操作部よりかな文字列の入力を受け付けた後に辞書検索に基づく変換処理を指示する操作を受け付けたことに応じて、入力されたかな文字列により変換用標準辞書および学習用辞書を検索して、変換後文字列の候補を抽出する候補検索手段と、前記操作部よりかな文字列の入力を受け付けた後に、カタカナ、ひらがな、アルファベット、数字のうちのいずれか一の文字種を選択して変換を指示する後変換操作を受け付けたことに応じて、入力されたかな文字列を指示された文字種による文字列に変換する後変換操作時処理手段と、この後変換操作時処理手段により変換され、前記変換用標準辞書および学習用辞書に登録されていない変換後文字列を学習用辞書に登録するとともに、前記形態素解析手段から検索対象の文字列を受け付けたとき、この文字列により学習用辞書を検索して、その検索により抽出した変換後文字列を形態素解析手段に渡す学習用辞書処理手段とを具備する。一方、前記形態素解析手段は、所定数の形態素が登録された形態素解析用辞書と、処理対象のテキストデータのを受け付けたとき、当該テキストデータから所定長さの文字列を抽出するステップAと、ステップAで抽出された文字列により形態素解析用辞書を検索して形態素の候補を抽出するステップBと、ステップAで抽出された文字列を文字変換手段の学習用辞書処理手段に渡し、当該文字列による学習用辞書の検索を実行した学習用辞書処理手段から渡された文字列を形態素の候補として受け付けるステップCとを複数サイクル実行する検索手段と、前記検索手段による各サイクルのステップBおよびステップCにより得た候補の中の所定数を選択して出力する出力手段とを具備する。

【0030】

さらに検索手段は、ステップAで抽出された文字列がひらがな文字列である場合には、ステップCにおいて、前記抽出された文字列を学習用辞書処理手段に渡さず、形態素の

10

20

30

40

50

候補を得られなかったという処理結果を設定する。また出力手段は、ステップCにより得たカタカナ文字列、アルファベット文字列、および数字による文字列を、ステップBにより得た候補より優先して選択する。

【0031】

上記の情報処理装置において、文字変換手段と形態素解析手段とは、いずれも、プログラムによって、情報処理装置の制御用コンピュータに設定されるものである。文字変換手段は、ワードプロセッサ、電子メールエディタなどのアプリケーションに連動して動作することができる。また、形態素解析手段は、ウェブブラウザやメールリーダーが動いているときに、これらからのコマンドやユーザーの呼び出し操作に応じて起動するものとして構成することができる。

10

【発明の効果】

【0034】

この発明によれば、ユーザの後変換操作に応じて自動的に学習された変換後文字列を、形態素解析処理用の辞書データとして使用することが可能になるから、形態素解析のために特別な登録作業を行ったり、形態素解析用辞書の更新データを取り込んだりしなくとも、形態素解析の機能を自然に向上することができる。

【0035】

また、この発明では、後変換操作に応じて学習用辞書に登録された変換後文字列のうちの重要な意味を持たない可能性が高いひらがな文字列が形態素の候補として抽出されないようにすることにより、学習用辞書から抽出される形態素の候補を、カタカナ文字列、アルファベット文字列、および数字による文字列に限定することができる。さらに学習用辞書から抽出された候補を形態素解析用辞書から抽出された候補より優先して選択するので、後変換操作により入力された特定の文字種（カタカナ、アルファベット、数字）による文字列を優先的に形態素の解析結果に含めることが可能になる。よって、たとえば、携帯電話において、受信メールを形態素解析し、その処理で抽出された単語を用いて返信メールを作成する際に、ユーザーが使用する可能性の高い単語を抽出するなど、ユーザーの関心に応じた新語を精度良く抽出することができる。また、システム開発者により形態素解析用辞書を更新するサービスを行わなくとも、各ユーザーは、自身の関心に応じた新語を抽出できるようになるなど、情報処理における利便性を大いに向上することができる。

20

【発明を実施するための最良の形態】

30

【0036】

図1は、この発明が適用された情報処理装置の構成例を示す。

この情報処理装置は携帯電話（図示せず。）に組み込まれるものであって、かな漢字変換処理部1と形態素解析処理部2とを含む。これらの処理部1, 2は、いずれも、プログラムによって、前記携帯電話の制御部（CPU）に設定されるものである。また、かな漢字変換処理部1は電子メールなどの文書作成用のアプリケーションとともに動作する。一方、形態素解析処理部2は、電子メールやウェブページの閲覧用のアプリケーションなどの上位アプリケーションが動いているときに、形態素解析処理を指定する操作に応じて起動する。

【0037】

40

かな漢字変換処理部1には、標準辞書16、ユーザー辞書17、自動学習辞書18の3種類の辞書や、入出力部11、変換制御部12が設けられるほか、各辞書毎に、その辞書に対する検索や登録のための処理部13, 14, 15が設けられる。いずれの辞書13, 14, 15にも、単語の表記（変換後文字列）に読みや使用頻度などを対応づけた辞書データが格納される。

【0038】

入出力部11は、前記携帯電話の表示部に展開されるヒューマンインターフェース（図示せず。）と情報をやりとりしつつ、操作キーの入力を受け付けるものである。たとえば、かな漢字変換処理のためにかな文字列（変換前文字列）が入力されている間は、操作されたキーの種類や操作回数に基づき、入力文字を認識し、その認識した文字をヒューマン

50

インターフェースに出力する。また、かな文字列の入力後に変換操作が行われると、その時点での未確定の変換後文字列をヒューマンインターフェースに出力する。さらに、この入出力部 11 は、未確定の変換後文字列が確定された場合には、文字入力対象のアプリケーション（メールエディタなど）に確定された文字列を出力する。

【0039】

この情報処理装置が導入される携帯電話の操作部では、前記変換操作や変換後文字列を確定するための操作のほか、後記する後変換操作のために、それぞれ特定のキーが設定される。前記入出力部 11 は、後変換用のキーが操作されると、前記ヒューマンインターフェースに文字種の選択画面を表示させ、この画面上で、カタカナ、ひらがな、アルファベット、数字のいずれかの文字種を選択させるようにしている。また、カタカナ、アルファベット、数字については、全角文字または半角文字を選択することもできる。

10

なお、以下で、「後変換操作」という場合には、前記後変換用のキーの操作から最終の文字種を選択する操作までを指すものとする。

【0040】

変換制御部 12 は、入出力部 11 を介してかな文字列の入力を受け付け、これを各辞書の処理部 13, 14, 15 に渡して検索処理を実行させる。そして、各処理部 13, 14, 15 から返された検索結果に基づき、前記かな文字列を所定の変換後文字列に変換し、その変換後文字列を入出力部 11 に出力する。

【0041】

3 種類の辞書のうち、標準辞書 16 には、システム設計者によりあらかじめ選択された単語の辞書データが格納される。この標準辞書 16 には、新たな辞書データを書き込むことはできないが、使用頻度については、辞書データの文字列が採用される都度、更新することができる。

20

【0042】

ユーザー辞書 17 は、ユーザーが指定した単語を登録するためのものである。変換制御部 12 は、登録対象の単語について、変換後文字列、読み、品詞情報などの入力を受け付けると、これをユーザー辞書処理部 14 に渡し、ユーザー辞書 17 に登録させる。

【0043】

自動学習辞書 15 は、漢字文字列以外の文字種への変換を指示する後変換操作が行われたときに、その操作による変換後文字列を登録するためのものである。この自動学習辞書 18 への登録処理も、変換制御部 12 により制御されるが、直接の登録処理は、自動学習辞書処理部 15 により行われる。

30

なお、ユーザー辞書 17、自動学習辞書 18 についても、登録されている辞書データが使用される都度、その使用頻度を更新することができる。

【0044】

つぎに、形態素解析処理部 2 には、入出力部 21、解析制御部 22、辞書検索部 23、形態素解析用辞書 24、候補制御部 25 などが含まれる。形態素解析用辞書 24 には、品詞情報、他の形態素との係り受けの関係などの属性データを形態素の表記に対応づけた辞書データが格納される。この形態素解析用辞書 24 も、前記かな漢字変換処理部 1 の標準辞書 16 と同様に、システム設計者により作成されたもので、新規の辞書データを登録するようには設定されていない。

40

【0045】

入出力部 21 は、前記した上位アプリケーションと連絡するもので、処理対象のテキストデータを入力し、そのテキストデータに対する形態素解析の結果（抽出した形態素を属性データとともに順に並べたもの）を出力するように、設定される。

【0046】

解析制御部 22 は、入出力部 21 から前記テキストデータの提供を受けて、このテキストデータから検索キーワードを切り出す。そして、この文字列と表記が一致する形態素を検索するために、辞書検索部 23、候補制御部 25、およびかな漢字変換処理部 1 のユーザー辞書処理部 14 に、前記検索キーワードを出力する。

50

【 0 0 4 7 】

辞書検索部 2 3 は、解析制御部 2 2 から検索キーワードを渡されると、形態素解析用辞書 2 4 を検索し、形態素の表記が前記検索キーワードに一致する辞書データを抽出する。ユーザー辞書処理部 1 4 も、同様に、解析制御部 2 2 から渡された検索キーワードを用いてユーザー辞書 1 7 を検索し、変換後文字列の表記が前記検索キーワードに一致する辞書データを抽出する。これらの検索で抽出された辞書データは、いずれも解析制御部 2 2 に渡される。

【 0 0 4 8 】

一方、候補制御部 2 5 は、かな漢字変換処理部 1 の自動学習辞書処理部 1 5 に前記検索キーワードを渡して、自動学習辞書 1 8 を用いた検索を実行させる。ただし、候補制御部 2 5 では、検索キーワードがひらがな文字列である場合には、これを自動学習辞書処理部 1 5 に渡さずに、解析制御部 2 2 に「候補なし」との検索結果を返す。ひらがな文字列以外の検索キーワードは、自動学習辞書処理部 1 5 に与えられるので、自動学習辞書処理部 1 5 は、与えられた検索キーワードにより自動学習辞書 1 8 を検索し、変換後文字列の表記が前記検索キーワードに一致する辞書データを抽出する。抽出された辞書データは、候補制御部 2 5 を介して解析制御部 2 2 に渡される。

なお、形態素解析処理における自動学習辞書 1 8 やユーザー辞書 1 7 の使用は、データの読み出しのみであって、データの書込みは行わないのが望ましい。

【 0 0 4 9 】

解析制御部 2 2 は、詳細は後記するが、複数の検索キーワードを設定し、これらの検索キーワード毎に、上記 3 通りの検索により形態素の候補を求める。さらに、解析制御部 2 2 は、検索キーワード毎に得た候補の中から、形態素の組み合わせとして最適なものを決定し、その決定に基づいて前記テキストデータの構造を表したデータ（以下、「解析データ」という。）を出力する。なお、解析データは、形態素の区切り位置のほか、各形態素の品詞情報などを含むものとなる。

【 0 0 5 0 】

以下、前記かな漢字変換処理部 1 および形態素解析処理部 2 で実行される処理のうち、この発明に関連する処理について、詳細に説明する。なお、図 2 以下のフローチャートでは、各ステップを「S T」と略記し、各ステップを 3 桁の番号で示す。

【 0 0 5 1 】

図 2 は、かな漢字変換処理部 1 による文字入力処理の手順を示す。なお、この実施例の手順では、複数文節分の文字列を入力して一括変換した後、最終的な確定操作が行われるまで、必要な文節を修正できるようにしている。

【 0 0 5 2 】

この手順は、かな文字の入力操作が開始されることによってスタートする。まず、最初の S T 1 0 1 では、ユーザーによる操作の内容を判断する。ユーザーは、文字入力のために、所定のかな文字が割り当てられたキーを操作するが、その操作の都度、S T 1 0 1 から S T 1 1 5 に進み、その操作に応じた処理（たとえば、入力文字をメモリ内の一時保存領域に蓄積する処理、蓄積された文字列を表示する処理など）を実行する。

【 0 0 5 3 】

所定数のかな文字が入力された時点で変換操作が行われると、S T 1 0 1 から S T 1 0 2 に進み、蓄積された入力文字列（かな文字列）を取得する。つぎの S T 1 0 3 では、この入力文字列を変換する処理（かな漢字変換処理）を実行する。このかな漢字変換処理は、公知の技術であるので、詳細は省略するが、入力文字列を所定数の文節に切り分け、標準辞書処理部 1 6、ユーザー辞書処理部 1 7、自動学習辞書処理部 1 8 に検索を実行させることによって、文節毎に最適な候補を抽出する。この実施例では、抽出された文節毎の候補を組み合わせた変換後文字列は、未確定の文字列として表示部に表示されるものとする。このとき、先頭の文節が注目文節として初期設定され、カーソルによる識別表示が行われる。

【 0 0 5 4 】

この実施例では、前記未確定の変換後文字列に含まれる文節毎に、後変換フラグを設定している。この後変換フラグは、後変換操作による変換処理を受けた旨を記憶するためのものであり、かな漢字変換処理の直後の S T 1 0 4 では、すべての文節の後変換処理をリセットする処理を実行する。

【 0 0 5 5 】

ここで変換後文字列がユーザーが目的とする文字列に一致していれば、確定操作が行われる。この操作により、S T 1 0 5 , 1 0 6 が「 Y E S 」となって S T 1 1 3 に進み、後記する自動学習処理を実行する。さらに、つぎの S T 1 1 4 では、前記未確定の変換後文字列を確定して出力する。これにより、入力文字列に対する処理は終了したことになる。

【 0 0 5 6 】

一方、前記未確定の変換後文字列の中に、ユーザーが意図しない文字列を含む文節がある場合には、ユーザーは、その文節にカーソルを移動させ、修正操作を実行する。これらの操作が行われる毎に、S T 1 0 5 が「 Y E S 」、S T 1 0 6 が「 N O 」となり、S T 1 0 7 以下の処理に進む。

【 0 0 5 7 】

まず、カーソルの移動操作が行われた場合には、S T 1 0 7 が「 Y E S 」となって S T 1 0 8 に進み、前記操作に応じて注目文節を変更する。このとき、前記カーソルも、移動操作に応じて注目文節に移動することになる。

【 0 0 5 8 】

つぎに、注目文節に対し、後変換操作が行われると、S T 1 0 9 が「 Y E S 」となって S T 1 1 0 に進み、注目文節を操作に応じた文字種に変更する。さらに、つぎの S T 1 1 1 では、前記注目文節の後変換フラグをオンにする。

【 0 0 5 9 】

また、注目文節に対し、後変換操作以外の操作（他の変換候補を選択する操作、変換をキャンセルする操作など）が行われた場合には、S T 1 1 2 に進み、その操作に応じた処理を実行する。修正処理が終了し、確定操作が行われると、前記した S T 1 1 3、1 1 4 を実行し、処理を終了する。

【 0 0 6 0 】

図 3 は、前記 S T 1 1 3 の自動学習処理について、詳細な手順を示す。

この処理では、S T 2 0 1 において、先頭の文節を注目文節に設定した後、S T 2 0 2 ~ 2 0 8 のループを繰り返すことにより、各文節を順に処理する。ただし、後変換フラグがオフの場合には、つぎの文節に注目文節を移すだけであり、後変換フラグがセットされている文節に対してのみ、実質的な処理（S T 2 0 3 ~ 2 0 6）を実行する。

【 0 0 6 1 】

注目文節に後変換フラグがセットされている場合、S T 2 0 3 では、その注目文節の読み（変換前のかな文字列）と変換後文字列とを対応づけた辞書データを設定する。つぎの S T 2 0 4 では、自動学習辞書 1 8 に空き領域があるか否かをチェックする。空き領域がある場合には、S T 2 0 5 をスキップして S T 2 0 6 に進み、その空き領域に新規の辞書データを登録する。なお、この登録の際に、辞書データには、使用頻度や登録順序が付加される。

【 0 0 6 2 】

自動学習辞書 1 8 に空き領域がない場合には、S T 2 0 4 から S T 2 0 5 に進み、前記した登録順序に基づき、辞書内の最も古い辞書データを抽出して、これを削除する。ただし、この辞書データの使用頻度が所定値以上である場合には、つぎに古い辞書データを抽出するものとする。これにより、登録後の経過時間が比較的長い辞書データの中から使用頻度が最も少ないものを削除することができる。この後は、S T 2 0 6 に進み、前記削除により空いた領域に新規の辞書データを登録する。

【 0 0 6 3 】

なお、自動学習辞書 1 8 内の辞書データは、登録順序や使用頻度に基づき、適宜ソートするようにしてもよい。

10

20

30

40

50

また、文字入力処理の手順は、上記図2, 3に限らず、たとえば、入力文字列を一括変換した後に、文節毎に確定できるようにしてもよい。この場合には、確定操作が行われる毎に、その操作にかかる文節の後変換フラグをチェックし、この後変換フラグがオンであれば、確定後の文字列を自動学習辞書18に登録する処理を行うことになる。

【0064】

図4は、前記自動学習辞書18のデータ構成例を示す。図中の「読み」は、前記した入力文字列に、「表記」は変換後文字列に、それぞれ対応する。読みには、ひらがな文字列のほか、半角数字による文字列も含まれている。表記には、後変換操作における選択に応じて、カタカナ、ひらがな、アルファベット、数字のいずれかの文字種による文字列が格納される。

10

【0065】

つぎに、図5を用いて、形態素解析処理の概要を説明する。

この実施例では、処理対象のテキストデータを文の単位に区切った上で、文毎に、形態素解析の一手法である最長一致法を実行するようにしている。すなわち、処理対象の文を構成する文字列(以下、「処理対象文字列」という。)内の各文字(図5の例の場合、「朝」「日」「が」「富」・・・の各文字)に順に注目して検索を実行する。この検索では、注目文字から文の末尾までの文字列を検索対象として、1番目の文字から所定位置の文字までの文字列に一致する文字列を抽出する。

【0066】

たとえば、図5の処理対象文字列中の4番目の文字「富」を例にして説明すると、この「富」から図示しない末尾の文字までの文字列を検索対象文字列として、前記「富」から所定位置の文字までの文字列に一致する文字列を抽出する検索を実行する。この結果、図5の例では、「富士山」が最も長い文字列として抽出され、以下、「富士」、「富」の2種類の文字列が抽出されている。

20

【0067】

処理対象文字列内の各文字毎に抽出された文字列は、それぞれ、その文字から始まる形態素の候補としてメモリに格納される。以下では、文字毎に抽出された候補が格納されるメモリ領域を「検索結果リスト」といい、すべての文字の検索結果リストを含むメモリ領域を「候補リスト」ということにする。

【0068】

このようにして候補リストが作成されると、処理対象文字列の文字の並びに沿って各検索結果リストをチェックし、所定の候補を形態素として選択する。この場合に、最長一致法では、文字列の長いものから優先して候補を選択し、先に選択した候補に含まれる文字について、候補の選択をスキップするようにしている。ただし、候補の選択は1組に限らず、複数とおりの候補の組み合わせを設定し、その中から形態素の区切りが最適なものを選択するようにしている。

30

【0069】

図6は、前記形態素解析処理部2における詳細な処理の手順を示す。なお、この手順は、テキストデータ中の1つの文に対するものである。テキストデータ中に複数の文が含まれる場合には、この図6の手順が文毎に実行されることになる。

40

【0070】

まず、最初のST301では、処理対象文字列の先頭の文字を注目文字として初期設定する。つぎのST302では、前記候補リストを設定するためのメモリ領域をクリアし、しかる後にST303に進む。

【0071】

ST303では、前記注目文字から処理対象文字列の末尾の文字までの文字列を切り出し、これを検索対象として設定する。続くST304では、この検索対象の文字列を用いた検索処理を実行する。この検索処理の詳細については後述する。

【0072】

検索が終了すると、ST305において、候補が抽出されたかどうかを判断する。この

50

判断結果が「NO」の場合には、ST306に進み、注目文字のみから成る文字列を候補として抽出する。検索処理で候補が抽出された場合には、このST306をスキップしてST307に進む。ST307では、抽出された候補を、注目文字の検索結果リストに格納する。ST306を実行した場合にも、ST307において、注目文字から成る文字列を検索結果リストに格納する。

【0073】

以下、ST308, 309によって注目文字を処理対象文字列の末尾の文字まで動かしながら、上記ST303~307の処理を繰り返す。この後は、ST308からST310に進み、前記した候補リストから候補の組み合わせを設定し、その中から最適なものを選択する処理を実行する。そして、ST311では、選択した候補の組み合わせに基づいて解析データを作成し、これを上位のアプリケーションに出力する。

10

【0074】

図7は、前記ST304の候補検索処理、すなわち、前記ST303で設定された検索対象の文字列を用いた検索処理の詳細な手順を示す。

まず、ST401では、現在の注目文字に対応する検索結果リストをクリアする。つぎに、ST402において、検索対象の文字列の長さをカウンタLに格納した後、ST403において、前記検索対象の文字列を検索キーワードとして設定する。

【0075】

以下では、設定された検索キーワードの末尾を一文字ずつ削除しながら、前記した各辞書に対する検索を実行する。まず、ST404では、候補制御部25に検索キーワードを出力することによって、自動学習辞書18を検索させ、つぎのST405で、その検索結果を取得する。なお、この検索結果に候補が含まれている場合には、ST405では、その候補を検索結果リストに格納する処理を実行する。

20

【0076】

ST406では、ユーザー辞書処理部14に検索キーワードを出力することによって、前記ユーザー辞書17を検索させ、つぎのST407で、その検索結果を取得する。このST407でも、検索結果に候補が含まれている場合には、その候補を検索結果リストに格納する処理を実行する。

【0077】

つぎのST408では、辞書検索部23に検索キーワードを出力することによって、前記形態素解析用辞書24を検索させ、続くST409で、その検索結果を取得する。このST409でも、検索結果に候補が含まれている場合には、その候補を検索結果リストに格納する処理を実行する。

30

【0078】

このようにして、各辞書に対する検索を終了すると、前記Lの値を1つ小さくする。更新後のLが0より大きければ、ST412に進み、検索キーワードの末尾文字を削除した後、ST404に戻る。Lの値が0となった場合は、処理を終了する。

【0079】

上記の形態素解析処理において、先に後変換操作によって確定され、自動学習辞書18に登録された文字列(以下、「自動登録文字列」という。)が処理対象文字列中に含まれている場合には、この自動登録文字列が検索キーワードとして設定されたときの前記ST404, 405の処理によって、自動登録文字列を検索結果リストに格納することができる。

40

【0080】

ただし、検索キーワードがひらがな文字列である場合には、前記候補制御部25は、この文字列による検索処理をキャンセルするから、この検索キーワードが自動登録文字列であっても、その文字列は候補として抽出されない。すなわち、自動学習辞書18から候補として抽出されるのは、カタカナ、アルファベット、数字の3種類の文字種による文字列に限定されることになる。また、重要な意味を持たない可能性が高いひらがな文字列について、自動学習辞書18による検索をスキップすることができるので、検索にかかる時間

50

を短縮することができる。

【0081】

なお、前記図6のST310で、各検索結果リストの候補を組み合わせる際には、自動学習辞書18から抽出された候補を最優先し、つぎにユーザー辞書17から抽出された候補を選択するようにしている。よって、文字入力処理においてユーザーが後変換操作により入力した単語や、ユーザーの登録処理によってユーザー辞書に登録された単語を、形態素解析処理でも優先的に抽出することができる。

【0082】

また、この実施例では、自動学習辞書18には、後変換操作により変換された文字列のみが登録されるものとしたが、これに限らず、標準辞書16やユーザー辞書17に登録されていない漢字文字列を登録することもできる。たとえば、1文節として変換された文字列が意図しない文字列であったために、ユーザーが入力文字列を2つに区切って文節毎に変換操作を行うことにより、所望の文字列を入力した場合、確定された漢字文字列を自動学習辞書18に登録することができる。このように自動学習辞書18に漢字文字列を登録した場合には、形態素解析処理でも、これらの漢字文字列まで抽出できるようにするとよい。

10

【0083】

また、文字変換処理では、自動学習辞書18に登録された文字列を通常の変換操作により呼び出して使用した場合に、使用頻度を更新することができる。よって、形態素解析処理でも、自動学習辞書18に対する検索によって自動登録文字列を抽出した場合には、その文字列の使用頻度を考慮して候補の組み合わせ処理を行うようにしてもよい。

20

【図面の簡単な説明】

【0084】

【図1】この発明が適用された情報処理装置の機能ブロック図である。

【図2】文字入力処理の手順を示すフローチャートである。

【図3】自動学習処理の手順を示すフローチャートである。

【図4】自動学習辞書のデータ構成を示す図である。

【図5】形態素解析処理の概要を説明する図である。

【図6】形態素解析処理の手順を示すフローチャートである。

【図7】候補検索処理の手順を示すフローチャートである。

30

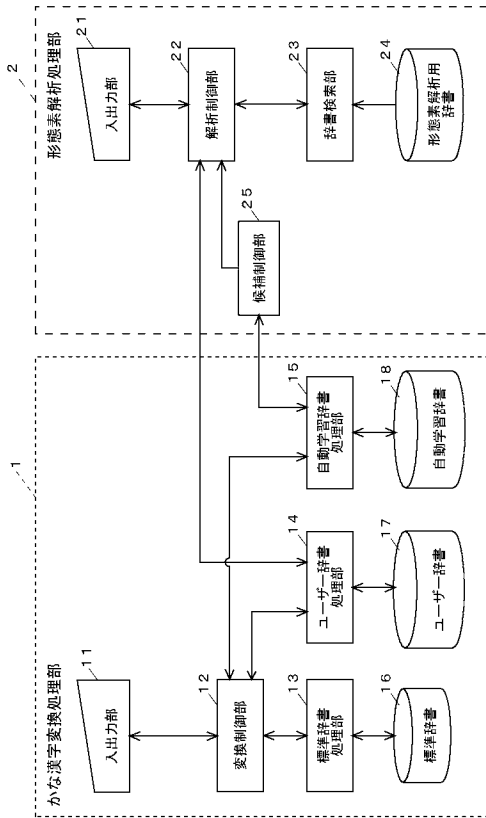
【符号の説明】

【0085】

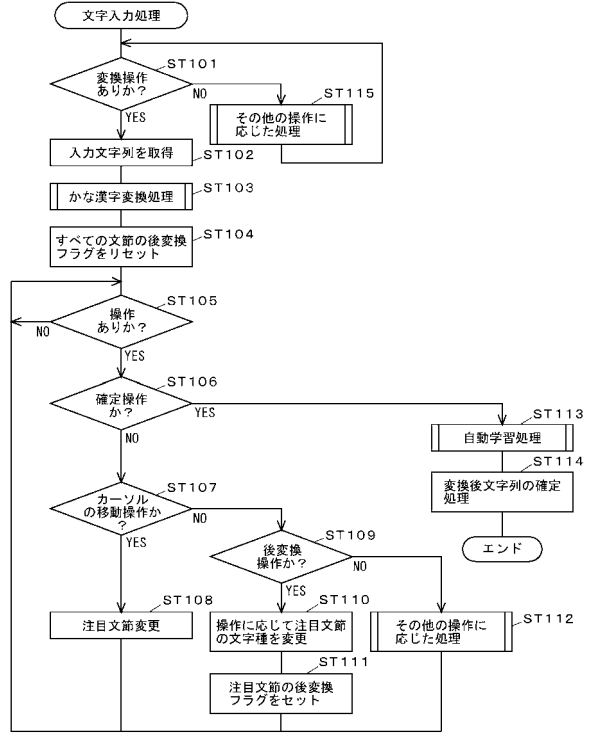
- 1 かな漢字変換処理部
- 2 形態素解析処理部
 - 1 1、2 1 入出力部
 - 1 2 変換制御部
 - 1 5 自動学習辞書処理部
 - 1 6 標準辞書
 - 1 8 自動学習辞書
- 2 2 解析制御部
- 2 3 辞書検索部
- 2 4 形態素解析用辞書
- 2 5 候補制御部

40

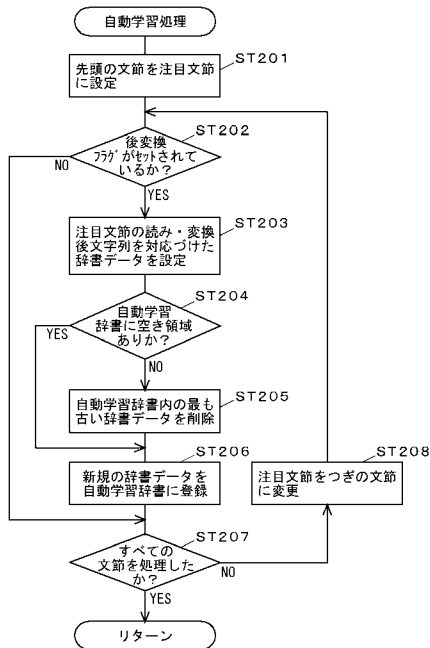
【図1】



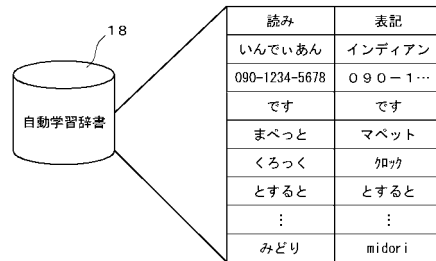
【図2】



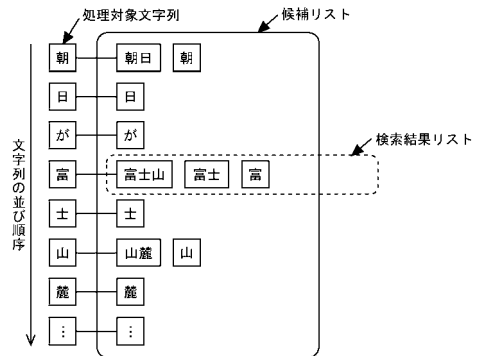
【図3】



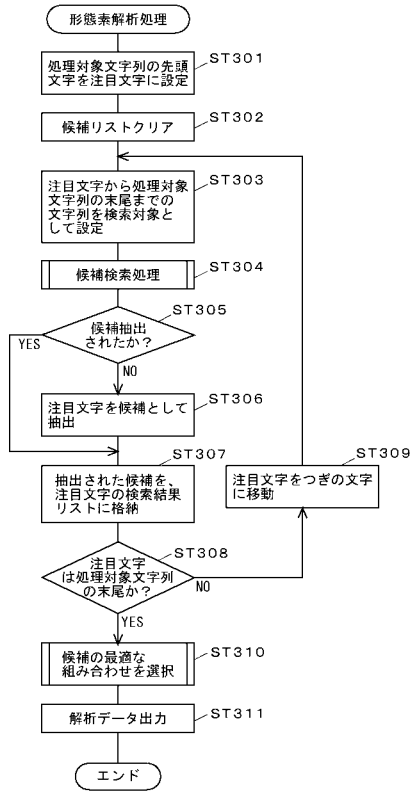
【図4】



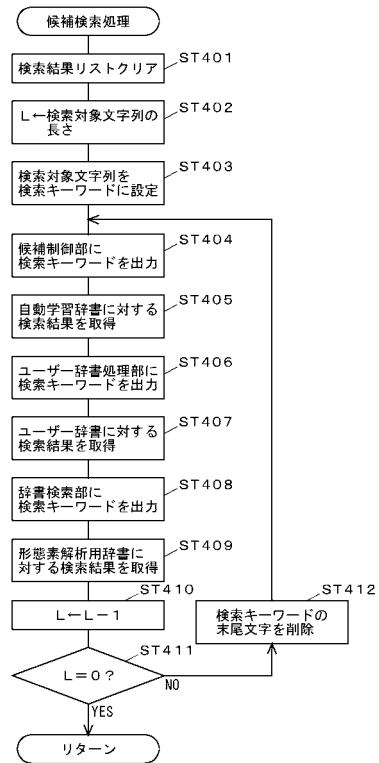
【図5】



【図6】



【図7】



フロントページの続き

審査官 今村 剛

(56)参考文献 特開2003-058537(JP,A)
特開平02-297151(JP,A)
特開平08-153098(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F17/20-17/28