

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5495425号  
(P5495425)

(45) 発行日 平成26年5月21日(2014.5.21)

(24) 登録日 平成26年3月14日(2014.3.14)

(51) Int.Cl.

F I

G06F 17/21 (2006.01)

G06F 17/24 (2006.01)

G06F 17/27 (2006.01)

G06F 17/30 (2006.01)

G06F 17/21 550A

G06F 17/24 554M

G06F 17/27 Z

G06F 17/30 170A

G06F 17/30 210A

請求項の数 11 (全 16 頁) 最終頁に続く

(21) 出願番号 特願2009-265784 (P2009-265784)  
 (22) 出願日 平成21年11月21日(2009.11.21)  
 (65) 公開番号 特開2011-113097 (P2011-113097A)  
 (43) 公開日 平成23年6月9日(2011.6.9)  
 審査請求日 平成24年9月4日(2012.9.4)

特許法第30条第1項適用 平成21年8月20日 発行の「第8回情報科学技術フォーラム 講演論文集(DVD)」に発表

(出願人による申告)平成21年度、独立行政法人情報通信研究機構「高度通信・放送研究開発委託研究／インターネット上の違法・有害情報の検出技術の研究開発」、産業技術力強化法第19条の適用を受ける特許出願

(73) 特許権者 599108264  
 株式会社KDDI研究所  
 埼玉県ふじみ野市大原二丁目1番15号  
 (74) 代理人 100135068  
 弁理士 早原 茂樹  
 (72) 発明者 池田 和史  
 埼玉県ふじみ野市大原二丁目1番15号  
 株式会社KDDI研究所内  
 (72) 発明者 柳原 正  
 埼玉県ふじみ野市大原二丁目1番15号  
 株式会社KDDI研究所内  
 (72) 発明者 松本 一則  
 埼玉県ふじみ野市大原二丁目1番15号  
 株式会社KDDI研究所内

最終頁に続く

(54) 【発明の名称】 未知語を含む文章を修正するための文章修正プログラム、方法及び文章解析サーバ

(57) 【特許請求の範囲】

【請求項1】

未知語を含む解析対象文章情報に対して、該未知語を修正するようにコンピュータを機能させる文章修正プログラムであって、

修正基準文章情報の集合を記憶した基準文章記憶手段と、

前記解析対象文章情報を形態素に分割し、未知語を抽出する未知語抽出手段と、

前記未知語に対応するワイルドカードと、当該未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる検索キーを生成する検索キー生成手段と、

前記基準文章記憶手段を用いて、前記解析対象文章情報の内容に類似する1つ以上の修正基準文章情報を検索する基準文章検索手段と、

前記修正基準文章情報の中から、前記検索キーを用いて、前記ワイルドカードで検索された1つ以上の修正ルールを検索する修正ルール検索手段と、

前記修正ルールの中から、文章構成指標に基づいて1つの修正ルールを選択する修正ルール選択手段と

してコンピュータを機能させることを特徴とする文章修正プログラム。

【請求項2】

前記基準文章記憶手段は、複数の修正基準文章情報を、その内容に基づいて複数のカテゴリに分類しており、

前記基準文章検索手段は、前記解析対象文章情報の内容に最も類似する前記カテゴリを検索し、当該カテゴリに含まれる前記修正基準文章情報を出力する

ようにコンピュータを機能させることを特徴とする請求項 1 に記載の文章修正プログラム。

【請求項 3】

前記基準文章検索手段は、前記解析対象文章情報における特徴語を抽出し、該特徴語をキーとして 1 つ以上の修正基準文章情報を検索するようにコンピュータを機能させることを特徴とする請求項 1 に記載の文章修正プログラム。

【請求項 4】

前記解析対象文章情報は、ネットワークを介して公開されているブログ(Weblog)、掲示板及び / 又はクチコミコメントにおける不特定多数のユーザによって記述された文章情報であるようにコンピュータを機能させることを特徴とする請求項 2 又は 3 に記載の文章修正プログラム。

10

【請求項 5】

前記基準文章記憶手段は、複数の修正基準文章情報を、その日時情報に基づいて複数のカテゴリに分類しており、

前記基準文章検索手段は、前記解析対象文章情報に記述された日時情報、又は、前記解析対象文章情報が作成された日時情報に最も近い前記カテゴリを検索し、当該カテゴリに含まれる前記修正基準文章情報を出力する

ようにコンピュータを機能させることを特徴とする請求項 1 に記載の文章修正プログラム。

【請求項 6】

20

前記解析対象文章情報は、ネットワークを介して公開されているブログ、掲示板及び / 又はクチコミコメントにおける不特定多数のユーザによって記述された文章情報であり、

前記解析対象文章情報が前記ブログである場合、前記日時情報は、当該ブログの URL (Uniform Resource Locator) アドレスに含まれたものである

ようにコンピュータを機能させることを特徴とする請求項 5 に記載の文章修正プログラム。

【請求項 7】

前記基準文章記憶手段は、日時情報に基づいて分類された前記カテゴリ毎に、更に、修正基準文章情報の内容に基づいて複数のカテゴリに分類しており、

前記基準文章検索手段は、日時情報によって検索された前記カテゴリの中から、前記解析対象文章情報の内容に最も類似する前記カテゴリを検索し、当該カテゴリに含まれる前記修正基準文章情報を出力する

30

ようにコンピュータを機能させることを特徴とする請求項 5 又は 6 に記載の文章修正プログラム。

【請求項 8】

前記基準文章検索手段は、日時情報によって検索された前記カテゴリの中から、前記解析対象文章情報における特徴語を抽出し、該特徴語をキーとして 1 つ以上の修正基準文章情報を検索するようにコンピュータを機能させることを特徴とする請求項 5 又は 6 に記載の文章修正プログラム。

【請求項 9】

40

前記修正ルール選択手段は、前記文章構成指標として、( 1 ) 当該修正ルールにおける出現頻度、( 2 ) 前記未知語と前記修正ルールに基づく修正形態素との間の編集距離、及び / 又は、( 3 ) 修正前と修正後との形態素解析コスト値の差分、に基づいて 1 つの修正ルールを選択するようにコンピュータを機能させることを特徴とする請求項 1 から 8 のいずれか 1 項に記載の文章修正プログラム。

【請求項 10】

未知語を含む解析対象文章情報に対して、該未知語を修正する装置の文章修正方法であって、

修正基準文章情報の集合を、基準文章記憶部に記憶しており、

前記解析対象文章情報を形態素に分割し、未知語を抽出する第 1 のステップと、

50

前記未知語に対応するワイルドカードと、当該未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる検索キーを生成する第2のステップと、

前記基準文章記憶部を用いて、前記解析対象文章情報の内容に類似する1つ以上の修正基準文章情報を検索する第3のステップと、

前記修正基準文章情報の中から、前記検索キーを用いて、前記ワイルドカードで検索された1つ以上の修正ルールを検索する第4のステップと、

前記修正ルールの中から、文章構成指標に基づいて1つの修正ルールを選択する第5のステップと

を有することを特徴とする装置の文章修正方法。

【請求項11】

未知語を含む解析対象文章情報を他の公開サーバからネットワークを介して取得し、該未知語を修正する文章解析サーバであって、

修正基準文章情報の集合を記憶した基準文章記憶手段と、

前記解析対象文章情報を形態素に分割し、未知語を抽出する未知語抽出手段と、

前記未知語に対応するワイルドカードと、当該未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる検索キーを生成する検索キー生成手段と、

前記基準文章記憶手段を用いて、前記解析対象文章情報の内容に類似する1つ以上の修正基準文章情報を検索する基準文章検索手段と、

前記修正基準文章情報の中から、前記検索キーを用いて、前記ワイルドカードで検索された1つ以上の修正ルールを検索する修正ルール検索手段と、

前記修正ルールの中から、文章構成指標に基づいて1つの修正ルールを選択する修正ルール選択手段と

を有することを特徴とする文章解析サーバ。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、未知語を含む文章を修正するための文章修正プログラム、方法及び文章解析サーバに関する。

【背景技術】

【0002】

インターネットの普及により、ブログ、掲示板又はクチコミコメントを公開するWebサイトに、様々なテキストが記述されている。「ブログ」(Weblog)とは、一般的に個人によって運営され、時事ニュースや専門的トピックスに関する自らの意見を表明するために、日記的に更新することができるサイトをいう。また、「掲示板」とは、様々なテーマについて、他人と議論を逐次に交換するためのサイトをいう。更に、「クチコミコメント」とは、人の噂のような、物事の評判などに関するコメントを記述することができるサイトをいう。これらサイトは、親しみやすさから、口語的な「くだけた表現」で記述されることが多い。

【0003】

近年、このようなサイトによって公開されるWeb文書も、情報抽出、検索及び統計処理の対象とされるようになってきた。これらWeb文書の文章内容を解析するために、少なくとも形態素解析が必要となる。形態素解析プログラムは、解析対象となる文章を形態素に分割する。英語の文章は、“This is a pen.”のように単語ごとに区切られた「分かち書き」にされている。これに対して、日本語の文章は、分かち書きされていないため、構文解析等に先立って、形態素解析による処理が必要となる。

【0004】

「形態素」とは、文章の構成要素のうち、意味を持つ最小の単位をいう。形態素解析プログラムは、「単語」毎に「品詞」「読み」が登録された辞書を有する。分割された形態素には、辞書を用いて「品詞」「読み」の情報が付与され、単語の配列を規定する文法規則を用いて文章を解析する。このように、形態素解析プログラムの解析精度は、辞書に依

10

20

30

40

50

存する。そのため、辞書に登録されていない単語を含む文章は、十分な解析精度を得ることができない。

【0005】

これに対し、ブログ、掲示板又はクチコミコメントによって公開されるWeb文書には、一般的に、以下のような「くだけた表現」が多数含まれる。

(1)「うっそー」「すごーい」のような会話における発音の変化傾向に併せた表記

(2)「カッコイイ」のように本来ひらがなで表記される語を意図的にカタカナにした表記

(3)「かわいい」(「かわいい」と読む)、「わたしわ」(「わたしは」と読む)のような特有の表記

10

【0006】

これら表現は、一般的に、形態素解析プログラムの辞書には登録されていないために、「未知語」として処理される。辞書に登録されていない単語は、形態素相当の単位に分割された上で、「品詞」「読み」の情報に代えて「未知語」という情報のみが付与される。

【0007】

勿論、「未知語」として処理された語を全て、人手によって辞書に登録することができれば、形態素解析の精度を向上させることができる。しかし、「未知語」の登録には、品詞及び活用形の登録、既存の辞書との互換性の維持といった、専門的な人手のスキルが必要となる。

【0008】

20

これに対して、未知語を含む文章を形態素解析に適した文章に修正する技術がある(例えば非特許文献1参照)。この技術によれば、解析に失敗して「未知語」として出力されることを前提として、未知語を含む文章である解析対象文章を形態素解析する。そして、形態素解析の結果に基づいて、未知語の部分を任意文字列に置き換えた検索キーを生成する。生成された検索キーによって、修正基準文章から、自動的に修正候補文字列が検索され、修正候補文字列から1以上の修正ルールが生成される。生成された修正ルールは、(1)同じような文脈で頻繁に使用される表現か、(2)解析対象文章の表現から変化し過ぎていないか、(3)修正後の文章が日本語として自然か、といった指標を用いて、スコアリングされ、最適な修正ルールが選択される。解析対象文章は、選択された修正ルールによって、形態素解析に適した文章に修正される。

30

【先行技術文献】

【非特許文献】

【0009】

【非特許文献1】池田和史、柳原正、松本一則、滝嶋康弘、「くだけた表現を修正するための教師なし学習方式の提案と評価」、第8回情報科学技術フォーラム、2009

【発明の概要】

【発明が解決しようとする課題】

【0010】

非特許文献1に記載された技術によれば、修正基準文章から検索された修正候補文字列に基づいて修正ルールが生成されるため、解析対象文章は、修正基準文章に依存して修正される。その結果、解析対象文章は、正しく形態素解析される文章には修正されても、異なる意味内容の文章に修正される場合や、過剰修正される場合があった。

40

【0011】

また、非特許文献1に記載された技術によれば、修正基準文章には、例えば、新聞記事の文章のような形態素解析の解析精度が高い文章が用いられている。一般に、新聞記事の文章は、「定型的な表現」が多い。「定型的な表現」は、形態素解析プログラムの辞書に登録されている確率が高いため、新聞記事の文章は、形態素解析の精度が高い。

【0012】

例えば、解析対象文章の「えーゆーはかっこいい」という表現は、新聞記事の文章を修正基準文章とすると、「英雄はかっこいい」と修正される。この解析対象文章がITやコ

50

ンピュータに関する文章であれば、この表現は、「a u はかっこいい」と修正されることが望ましい。

【0013】

また、例えば、解析対象文章の「そんなの関係ねえ」という表現は、「そんなの関係ない」と修正される。しかし、この解析対象文章が数年前の流行を反映した文章であれば、この表現は、「そんなの関係ねえ」のまま修正されないことが望ましい。

【0014】

そこで、本発明は、口語的な「くだけた表現」で記述された文章情報であっても、形態素解析によって未知語と判断されることのない、文章解析に適した文章に修正することができる文章修正プログラム、方法及び文章解析サーバを提供することを目的とする。

10

【課題を解決するための手段】

【0015】

本発明によれば、未知語を含む解析対象文章情報に対して、該未知語を修正するようにコンピュータを機能させる文章修正プログラムであって、

修正基準文章情報の集合を記憶した基準文章記憶手段と、

解析対象文章情報を形態素に分割し、未知語を抽出する未知語抽出手段と、

未知語に対応するワイルドカードと、当該未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる検索キーを生成する検索キー生成手段と、

基準文章記憶手段を用いて、解析対象文章情報の内容に類似する1つ以上の修正基準文章情報を検索する基準文章検索手段と、

20

修正基準文章情報の中から、検索キーを用いて、ワイルドカードで検索された1つ以上の修正ルールを検索する修正ルール検索手段と、

修正ルールの中から、文章構成指標に基づいて1つの修正ルールを選択する修正ルール選択手段と

してコンピュータを機能させることを特徴とする。

【0016】

本発明の文章修正プログラムにおける他の実施形態によれば、

基準文章記憶手段は、複数の修正基準文章情報を、その内容に基づいて複数のカテゴリに分類しており、

基準文章検索手段は、解析対象文章情報の内容に最も類似するカテゴリを検索し、当該カテゴリに含まれる修正基準文章情報を出力するようにコンピュータを機能させることも好ましい。

30

【0017】

本発明の文章修正プログラムにおける他の実施形態によれば、基準文章検索手段は、解析対象文章情報における特徴語を抽出し、該特徴語をキーとして1つ以上の修正基準文章情報を検索するようにコンピュータを機能させることも好ましい。

【0018】

本発明の文章修正プログラムにおける他の実施形態によれば、解析対象文章情報は、ネットワークを介して公開されているブログ(Weblog)、掲示板及び/又はクチコミコメントにおける不特定多数のユーザによって記述された文章情報であるようにコンピュータを機能させることも好ましい。

40

【0019】

本発明の文章修正プログラムにおける他の実施形態によれば、基準文章記憶手段は、複数の修正基準文章情報を、その日時情報に基づいて複数のカテゴリに分類しており、

基準文章検索手段は、解析対象文章情報に記述された日時情報、又は、解析対象文章情報が作成された日時情報に最も近いカテゴリを検索し、当該カテゴリに含まれる修正基準文章情報を出力するようにコンピュータを機能させることも好ましい。

【0020】

本発明の文章修正プログラムにおける他の実施形態によれば、

解析対象文章情報は、ネットワークを介して公開されているブログ、掲示板及び/又は

50

クチコミコメントにおける不特定多数のユーザによって記述された文章情報であり、

解析対象文章情報がブログである場合、日時情報は、当該ブログのURL (Uniform Resource Locator) アドレスに含まれたものであるようにコンピュータを機能させることも好ましい。

【0021】

本発明の文章修正プログラムにおける他の実施形態によれば、

基準文章記憶手段は、日時情報に基づいて分類されたカテゴリ毎に、更に、修正基準文章情報の内容に基づいて複数のカテゴリに分類しており、

基準文章検索手段は、日時情報によって検索されたカテゴリの中から、解析対象文章情報の内容に最も類似するカテゴリを検索し、当該カテゴリに含まれる修正基準文章情報を出力するようにコンピュータを機能させることも好ましい。

10

【0022】

本発明の文章修正プログラムにおける他の実施形態によれば、基準文章検索手段は、日時情報によって検索されたカテゴリの中から、解析対象文章情報における特徴語を抽出し、該特徴語をキーとして1つ以上の修正基準文章情報を検索するようにコンピュータを機能させることも好ましい。

【0023】

本発明の文章修正プログラムにおける他の実施形態によれば、修正ルール選択手段は、文章構成指標として、(1) 当該修正ルールにおける出現頻度、(2) 未知語と修正ルールに基づく修正形態素との間の編集距離、及び/又は、(3) 修正前と修正後との形態素解析コスト値の差分、に基づいて1つの修正ルールを選択するようにコンピュータを機能させることも好ましい。

20

【0024】

本発明によれば、未知語を含む解析対象文章情報に対して、該未知語を修正する装置の文章修正方法であって、

修正基準文章情報の集合を、基準文章記憶部に記憶しており、

解析対象文章情報を形態素に分割し、未知語を抽出する第1のステップと、

未知語に対応するワイルドカードと、当該未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる検索キーを生成する第2のステップと、

基準文章記憶部を用いて、解析対象文章情報の内容に類似する1つ以上の修正基準文章情報を検索する第3のステップと、

30

修正基準文章情報の中から、検索キーを用いて、ワイルドカードで検索された1つ以上の修正ルールを検索する第4のステップと、

修正ルールの中から、文章構成指標に基づいて1つの修正ルールを選択する第5のステップと

を有することを特徴とする。

【0025】

本発明によれば、未知語を含む解析対象文章情報を他の公開サーバからネットワークを介して取得し、該未知語を修正する文章解析サーバであって、

修正基準文章情報の集合を記憶した基準文章記憶手段と、

40

解析対象文章情報を形態素に分割し、未知語を抽出する未知語抽出手段と、

未知語に対応するワイルドカードと、当該未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる検索キーを生成する検索キー生成手段と、

基準文章記憶手段を用いて、解析対象文章情報の内容に類似する1つ以上の修正基準文章情報を検索する基準文章検索手段と、

修正基準文章情報の中から、検索キーを用いて、ワイルドカードで検索された1つ以上の修正ルールを検索する修正ルール検索手段と、

修正ルールの中から、文章構成指標に基づいて1つの修正ルールを選択する修正ルール選択手段と

を有することを特徴とする。

50

## 【発明の効果】

## 【0026】

本発明の文章修正プログラム、方法及び文章解析サーバによれば、口語的な「くだけた表現」で記述された文章情報であっても、形態素解析によって未知語と判断されることのない、文章解析に適した文章に修正することができる。

## 【図面の簡単な説明】

## 【0027】

【図1】本発明における文章修正プログラムの機能構成図である。

【図2】本発明におけるカテゴリに基づいて基準文章を検索する説明図である。

10

【図3】本発明における特徴語に基づいて基準文章を検索する説明図である。

【図4】本発明における日時情報に基づいて基準文章を検索する説明図である。

【図5】本発明における文章解析サーバのシステム構成図である。

【図6】本発明におけるシステムのシーケンス図である。

## 【発明を実施するための形態】

## 【0028】

以下、本発明の実施の形態について、図面を用いて詳細に説明する。

## 【0029】

図1は、本発明における文章修正プログラムの機能構成図である。

## 【0030】

20

図1によれば、文章修正プログラム1は、基準文章記憶部11と、基準文章検索部12と、未知語抽出部13と、検索キー生成部14と、修正ルール検索部15と、修正ルール選択部16と、修正ルール適用部17とを有する。基準文章記憶部11を除くこれら機能部は、装置に搭載されたコンピュータを機能させるプログラムを実行することによって実現できる。尚、各機能部の処理の流れは、コンピュータを用いた文章修正方法として実行できる。

## 【0031】

基準文章記憶部11は、修正基準文章情報の集合を記憶する。修正基準文章は、例えば技術文書、ブログテキスト、雑誌記事及び新聞記事のような様々な分野の文章を含む。基準文章記憶部11は、修正基準文章情報を、文章内容のカテゴリ又は特徴語によって分類していてもよいし、文章の内容的日時（又は作成日時）によって分類していてもよい。基準文章記憶部11は、基準文章検索部12によって参照される。

30

## 【0032】

未知語抽出部13は、解析対象文章情報を入力する。解析対象文章情報は、Webサイトに公開されているブログ、掲示板又はクチコミコメントのような不特定多数のユーザによって記述された文章情報であってもよい。未知語抽出部13は、その解析対象文章情報を形態素解析によって形態素に分割する。ここで、「できるかどうか分かりません」というくだけた表現を含む解析対象文章を例に挙げて説明する。

解析対象文章：できるかどうか分かりません

形態素解析結果：できる / か / どう / カわ（未知語） / 分かり / ませ / ん

40

くだけた表現は、形態素解析辞書に登録されていない場合が多い。そこで、形態素解析辞書に登録されていない表現「カわ」は、未知語として処理される。

## 【0033】

未知語抽出部13は、解析対象文章から未知語を検出した場合、形態素解析によって抽出された未知語と、未知語に隣接する形態素とを合わせた文字列とを、検索キー生成部14へ出力する。また、未知語抽出部13は、未知語が検出された解析対象文章情報を、基準文章検索部12へ出力する。更に、未知語抽出部13は、未知語が検出された解析対象文章情報を修正ルール選択部16と、修正ルール適用部17とへ出力する。

## 【0034】

基準文章検索部12は、解析対象文章情報を入力する。また、基準文章検索部12は、

50

基準文章記憶部 11 を参照し、解析対象文章の内容に類似する 1 つ以上の修正基準文章情報を検索によって選択する。基準文章検索部 12 は、基準文章記憶部 11 の分類に基づいて、例えば以下の 3 つの条件で検索する。

(1) 解析対象文章情報の内容に最も類似するカテゴリに対応した修正基準文章情報を検索する。

(2) 解析対象文章情報における特徴語を抽出し、その特徴語をキーとして修正基準文章情報を検索する。

(3) 解析対象文章情報に記述された日時情報、又は、解析対象文章情報が作成された日時情報に最も近い修正基準文章情報を検索する。

そして、基準文章検索部 12 は、1 つ以上の修正基準文章情報を、修正ルール検索部 15 へ出力する。

#### 【0035】

検索キー生成部 14 は、未知語と、当該未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる「検索キー」を生成する。ここで、未知語と、それに隣接する前方後方の各 1 形態素とが、検索キーの生成に利用されたものとして説明する。

入力文字列 : どう / カわ (未知語) / 分かり

検索キー : どう \* 分かり (ここで、「\*」は、1 以上の任意文字列を示す。)

#### 【0036】

検索キー生成部 14 は、未知語を任意文字列 (例えばワイルドカード) とし、任意文字列と、未知語に隣接する文字列と合わせた検索キーを生成する。勿論、未知語は、2 以上連続するものであってもよい。また、隣接する形態素は、未知語に対する前方及び後方の少なくとも一方があればよい。同様に、隣接する形態素も、2 形態素以上連続するものであってもよい。検索キー生成部 14 は、生成した検索キーを修正ルール検索部 15 へ出力する。

#### 【0037】

修正ルール検索部 15 は、検索キーと 1 つ以上の修正基準文章情報とを入力する。修正ルールとは、未知語 (例えば「かわいい」) から、修正候補文字列 (例えば、「かわいい」) へ文字列変換するためのルールをいう。修正ルール検索部 15 は、修正基準文章情報の中から、検索キーを含む修正候補文字列を検索する。そして、修正ルール検索部 15 は、抽出した修正候補文字列中の任意文字列に該当する部分を、未知語に近似する部分と判断し、修正ルールとして抽出する。

#### 【0038】

修正ルール検索部 15 は、例えば、検索により、以下の修正候補文字列を得る。修正ルール検索部 15 は、抽出した修正候補文字列中の任意文字列に該当する部分から、1 以上の修正ルールを抽出する。

検索キー : どう \* 分かり (ここで、「\*」は、1 以上の任意文字列を示す。)

修正候補文字列 : これはどう / かは / 分かりません

よくあるかどうか / か / 分かりません

どう / したらいいのか / 分かりません

この先どう / かは / 分かりません

本当かどうか / か / 分かりませんが

使うかどうか / かは / 分かりませんがね

あるかどうか / かは / 分かりません

どう / なっているか / 分かりませんよ

修正ルール : カわ かは

かわ か

かわ したらいいのか

かわ なっているか

#### 【0039】

検索キーによる検索によって得られる修正ルールは、2 以上であってもよい。修正ル

10

20

30

40

50



ル検索部 15 は、検索キーによる検索によって得た全ての修正ルールを、修正ルール選択部 16 へ出力する。

【 0 0 4 0 】

修正ルール選択部 16 は、入力された修正ルールが 2 以上ある場合は、文章構成指標に基づいて文脈に適した 1 つの修正ルールを選択する。

【 0 0 4 1 】

文章構成指標は、( 1 ) 修正ルールにおける出現頻度、( 2 ) 未知語と、修正ルールに基づく修正形態素との間の編集距離、及び / 又は、( 3 ) 修正前と修正後との形態素解析コスト値の差分から算出される指標をいう。修正ルール選択部 16 は、この文章構成指標に基づいて 1 つの修正ルールを選択する。

10

【 0 0 4 2 】

修正ルール選択部 16 は、例えば、以下の修正ルールを入力したとする。

修正ルール      : カわ    かは  
                      かわ    か  
                      かわ    したらいいのか  
                      かわ    なっているか

【 0 0 4 3 】

( 1 ) 修正ルールにおける出現頻度

修正ルールにおける出現頻度は、検索された修正ルールに該当する検索結果文字列が出現した頻度をいう。以下の表では、検索結果文字列の出現頻度に基づくスコアリングの例を表す。

20

【 表 1 】

修正候補文字列の出現頻度に基づくスコアリング例

修正ルール	出現頻度	出現頻度/検索件数
かわ⇒ かは	4	0.5
かわ⇒ か	2	0.25
かわ⇒ したらいいのか	1	0.125
かわ⇒ なっているか	1	0.125

【 0 0 4 4 】

30

出現頻度が高い文字列は、未知語が出現した文脈と類似した文脈の中で頻繁に利用される表現であると考えられ、修正候補文字列である可能性が高い。一方、類似した文脈の中であまり利用されていない表現は、修正候補文字列ではない可能性が高い。そこで、出現頻度の高い修正ルールは、スコアが高くなる。スコアは、出現頻度を検索件数で割り、正規化することにより、検索件数に依存しないものとしてもよい。

【 0 0 4 5 】

( 2 ) 未知語と、修正ルールに基づく修正形態素との間の編集距離

編集距離とは、二つの文字列がどの程度異なっているかを表す指標であり、一方の文字列を他方の文字列に変換するために必要な挿入、削除、置換の最小回数として与えられる。修正ルールに基づく修正形態素は、未知語に対して少数文字の挿入や削除、置換を実行したものであることが多い。例えば、「フォーラム」から「ファーム」への編集は、「ォ」を「ァ」に置換し、「ラ」を削除する方法が、最小の編集回数である 2 回となるため、編集距離は 2 である。以下の表は、編集距離に基づくスコアリングの例を表す。

40

【表 2】

編集距離に基づくスコアリング例

修正ルール	編集手順	編集距離
かわ⇒かは	置換:2回	2
かわ⇒か	置換:1回、削除:1回	2
かわ⇒したらいいのか	置換:2回、挿入:5回	7
かわ⇒なっているか	置換:2回、挿入:4回	6

## 【0046】

編集距離の小さい修正ルールは、スコアが高くなる。また、Web文書では、「ヤバイ」や「カッコイイ」のように本来ひらがなで表記されるべき語がカタカナで表記されている例が多い。そのため、例えば、カタカナをひらがなに置換する編集距離を小さくする重み付き編集距離を用いてもよい。

10

## 【0047】

## (3) 形態素解析コスト値の差分

形態素解析コスト値とは、複数ある単語区切りの中で、その単語区切りがどのくらい確からしいかを表す指標である。形態素解析コスト値は、例えば、単語単体での出現確率(生起コスト)や複数単語が連続して出現する確率(接続コスト)から算出される。形態素解析コスト値は、修正ルールの文脈における適応度を評価する指標として用いられる。

## 【0048】

20

文全体の形態素解析コスト値は、文頭から文末までの各形態素の接続コストと単語生起コストとの和を累積して算出する(累積コスト)。修正ルールの適用により、文脈における適応度が高い表現が生成された場合、その表現周辺の生起コストや接続コストは小さくなるため、文全体の形態素解析コスト値は小さくなる。一方、文脈における適応度が低い表現が生成された場合、その表現周辺の生起コストや接続コストは大きくなるため、文全体の形態素解析コスト値は大きくなる。

## 【0049】

ここでは、修正ルール適用後の文全体の形態素解析コスト値と、修正前の文全体の形態素解析コスト値との差分から、形態素解析コスト値に基づくスコアとして算出する。修正ルールによって生成された表現が文脈に適応する場合、算出されるスコアは高くなる。

30

## 【0050】

文章構成指標(score)は、(1)修正ルールにおける出現頻度(freq)、(2)未知語と、修正ルールに基づく修正形態素との間の編集距離(dist)、及び/又は、(3)形態素解析コスト値の差分(cost)から、例えば、以下の計算式により算出する。

$$\text{score} = \quad \cdot \text{freq} + \quad \cdot \text{dist} + \quad \cdot \text{cost}$$

ここで、 $\quad$ 、 $\quad$ 、 $\quad$ は、重み付け関数であり、修正ルールの適用と学習により、最適値を算出することができる。また、修正ルールは、適用する閾値を設定することができる。閾値を低く設定した場合、適用される修正ルールは増加するが、その中に含まれる修正ルールの誤適用も増加する。一方、閾値を高く設定した場合、適用される修正ルールは減少するが、その中に含まれる修正ルールの誤適用も減少させることができる。

40

## 【0051】

修正ルール選択部16は、文章構成指標に基づいて、1つの修正ルールを選択する。ここで、修正ルール選択部16は、閾値以上且つ最大のスコアを持つ修正ルールを選択してもよい。修正ルール選択部16は、選択した1つの修正ルールを、修正ルール適用部17へ出力する。

## 【0052】

修正ルール適用部17は、解析対象文章に修正ルールを適用する。修正ルール適用部17から出力された修正済みの文章情報は、様々な文章解析に適するものとなる。

## 【0053】

本発明の特徴は、解析対象文章に対して、適切な修正基準文章を検索によって選択する

50

ことにある。従って、基準文章記憶部 11 及び基準文章検索部 12 における複数の実施形態を、以下の図 2 ～ 図 4 によって説明する。

【0054】

図 2 は、本発明におけるカテゴリに基づいて基準文章を検索する説明図である。

【0055】

カテゴリに基づく基準文章検索では、予め複数の修正基準文章情報を、その内容に基づいて複数のカテゴリに分類している。基準文章のカテゴリは、文章中出现する語の偏りに基づいて分類される。

例えば、「着信履歴」「通話」といった語が、他の語と比較して多く出現する文章は、「携帯電話」のカテゴリに分類する。

10

また、「ウイルス」「スパイウェア」「ファイアウォール」といった語が、他の語と比較して多く出現する文章は、「情報セキュリティ」のカテゴリに分類する。

【0056】

カテゴリに基づく基準文章検索は、そのカテゴリに偏って多く出現した語をキーワードとして、解析対象文章をフィルタリングし、キーワードと一致する語数を計測する。ここで、解析対象文章は、キーワードと一致する語数が最も多いカテゴリに属する文章と推定される。基準文章検索部 12 は、解析対象文章情報の内容に最も類似するカテゴリを検索し、カテゴリに含まれる修正基準文章情報を出力する。

【0057】

図 3 は、本発明における特徴語に基づいて基準文章を検索する説明図である。

20

【0058】

特徴語に基づく基準文章検索では、解析対象文章から特徴語を抽出する。抽出された特徴語をキーワードとして検索することによって、関連性の高い文章を収集する。基準文章検索部 12 は、解析対象文章と関連性の高い 1 つ以上の修正基準文章情報を出力する。

【0059】

特徴語の抽出には、例えば  $tf \cdot idf$  法を用いることができる。 $tf \cdot idf$  法とは、文章中出现する単語の重み (= 特徴度合い) を計算する方法である。 $tf \cdot idf$  法は、文章中の特徴的な語を抽出するためのアルゴリズムとして用いられ、 $tf$  (Term Frequency、単語の出現頻度) と  $idf$  (Inverse Document Frequency、逆出現頻度) の 2 つの指標によって計算される。

30

【0060】

この方法によれば、単語  $t$  の文書  $d$  における重み  $w(t, d)$  は、次のように計算される。

$$w(t, d) = tf(t, d) \cdot idf(t)$$

$tf$  は、その単語が一つの文章中出现する頻度である。 $tf(d, t)$  は、文書  $d$  における単語  $t$  が現れる頻度を、文書内の形態素数で割った値である。 $tf$  は、文章中出现頻度が高い単語は、その文章において重要であると判断する指標である。

【0061】

一方、多くの文章中出现する単語は、文章を特定する性質を持たないことが多い。 $idf$  は、多くの文章中出现する語の重要度を下げ、その文章にのみ出現する単語の重要度を上げる指標である。 $idf$  は、文書の数  $N$  と、単語  $t$  が一回以上出現する文書の数によって、

40

以下の式のように定義される。

$$idf(t) = \log(N / df(t)) + 1$$

【0062】

図 4 は、本発明における日時情報に基づいて基準文章を検索する説明図である。

【0063】

日時情報に基づく基準文章検索では、予め複数の修正基準文章情報を、その日時情報に基づいて複数のカテゴリに分類している。基準文章のカテゴリは、文章の日時情報に基づいて分類される。文章の日時情報には、文章に記述された日時情報と、文章が作成された日時情報と、文章が更新された日時情報とがある。

【0064】

50

文章に記述された日時情報は、文章中に記載されている情報から判断される。例えば、「池田和史が、2003/12/09に、チーム松本に入団」と記載されている場合、文章に記述された日時情報は、2003年12月9日と判断される。

【 0 0 6 5 】

文章が作成された日時情報は、本文中に含まれている時間情報から判断される。例えば、「<http://www.blog.jp/20040105/index.html>」のようにURL (Uniform Resource Locator) アドレスに含まれていることが多い。この場合、文章が作成された日時情報は、2004年1月5日と判断される。

【 0 0 6 6 】

文章が更新された日時情報は、Webページのヘッダに含まれている時間情報から判断される。例えば、「Last-Modified: Tue, 19 Aug 2003 06:10:54 GMT」のような情報がWebページのヘッダに含まれている。この場合、文章が更新された日時情報は、2003年8月19日、6時10分54秒と判断される。

【 0 0 6 7 】

基準文章検索部 12 は、解析対象文章情報に記述された日時情報、又は、解析対象文章情報が作成された日時情報に最も近いカテゴリを検索し、そのカテゴリに含まれる修正基準文章情報を出力する。

【 0 0 6 8 】

また、基準文章検索部 12 は、解析対象文章情報が更新された日時情報に最も近いカテゴリを検索し、そのカテゴリに含まれる修正基準文章情報を出力してもよい。更に、日時情報として、文章に記述された日時情報、文章が作成された日時情報及び文章が更新された日時情報のうち複数が見られる場合、例えば、優先度をつけて、「文章に記述された日時情報」 > 「文章が作成された日時情報」 > 「文章が更新された日時情報」の順に判断してもよい。

【 0 0 6 9 】

基準文章検索部 12 は、予め複数の修正基準文章情報を、その日時情報に基づいて複数のカテゴリに分類したものを、更に、その修正基準文章の内容に基づいて複数のカテゴリに分類したものに対して、検索が実行されてもよい。また、基準文章検索部 12 は、予め複数の修正基準文章情報を、その日時情報に基づいて複数のカテゴリに分類したものに対して、解析対象文章から抽出した特徴語をキーワードとして、検索が実行されてもよい。

【 0 0 7 0 】

図 5 は、本発明における文章解析サーバのシステム構成図である。

【 0 0 7 1 】

図 5 によれば、文章解析サーバ 2 は、通信インタフェース部 20 と、基準文章入力部 21 と、解析対象文章入力部 22 と、文章修正機能部 23 と、基準文章記憶部 11 とを有する。文章解析サーバ 2 は、通信インタフェース部 20 を介してインターネットに接続する。

【 0 0 7 2 】

また、図 5 によれば、文章解析サーバ 2 は、インターネットを介して、Webサーバ 3 と通信することができる。また、Webサーバ 3 は、投稿者用端末 4 から接続される。

【 0 0 7 3 】

Webサーバ A は、ブログテキスト、雑誌記事、新聞記事又は技術文書のような様々な文章情報を公開している。文章解析サーバ 2 は、インターネットを介して、Webサーバ A から、それら文章情報を修正基準文章として収集する。

【 0 0 7 4 】

Webサーバ B は、投稿者用端末 4 から受信した、解析対象文章であるブログテキストを公開する。文章解析サーバ 2 は、インターネットを介して、Webサーバ B から、そのブログテキストを解析対象文章として取得する。

【 0 0 7 5 】

基準文章入力部 21 は、通信インタフェース部 20 を介して、修正基準文章を受信する

10

20

30

40

50

。その修正基準文章は、基準文章記憶部 11 へ出力される。

【0076】

解析対象文章入力部 22 は、通信インタフェース部 20 を介して、解析対象文章を受信する。その解析対象文章を、文章修正機能部 23 へ出力される。

【0077】

文章修正機能部 23 は、図 1 で前述した機能構成部と全く同様である。文章修正機能部 23 は、解析対象文章入力部 22 から解析対象文章を入力し、修正後文章を出力する。

【0078】

図 6 は、本発明におけるシステムのシーケンス図である。

【0079】

(S61) 文章解析サーバ 2 は、Webサーバ A から修正基準文章を収集する。それら修正基準文章は、基準文章記憶部 11 によって記憶される。

(S62) 投稿者用端末 4 は、解析対象文章であるブログテキストを Webサーバ B へ投稿する。

(S63) 文章解析サーバ 2 は、Webサーバ B から解析対象文章(「えーゆーはかっこいい」)を受信する。その解析対象文章は、文章修正機能部 23 へ出力される。

(S64) 文章修正機能部 23 は、解析対象文章情報を形態素に分割する。

(S65) 文章修正機能部 23 は、形態素に未知語が含まれていた場合、未知語を抽出する。また、文章修正機能部 23 は、未知語が検出された解析対象文章情報を、基準文章検索部 12 へ出力する。

(S66) 基準文章検索部 12 は、解析対象文章から、文章のカテゴリ、特徴語又は日時情報を抽出する。基準文章検索部 12 は、解析対象文章情報の内容に類似する 1 つ以上の修正基準文章情報を基準文章記憶部 11 から検索する。基準文章検索部 12 は、文章修正機能部 23 へ、修正基準文章情報を出力する。

(S67) 文章修正機能部 23 は、S65 で抽出した未知語と、その未知語に対する前方及び後方の少なくとも一方の隣接形態素とからなる検索キーを生成する。

(S68) 文章修正機能部 23 は、修正基準文章情報の中から、検索キーを用いて、未知語に近似する 1 つ以上の修正ルールを検索する。

(S69) 文章修正機能部 23 は、S68 で検索した修正ルールの中から、文章構成指標に基づいて 1 つの修正ルールを選択する。文章修正機能部 23 は、選択した修正ルールを解析対象文章に適用し、修正後文章(「au はかっこいい」)を出力する。

【0080】

以上、詳細に説明したように、本発明の文章修正プログラム、方法及び文章解析サーバによれば、口語的な「くだけた表現」で記述された文章情報であっても、形態素解析によって未知語と判断されることのない、文章解析に適した文章に修正することができる。

【0081】

解析対象文章に類似するカテゴリの文章集合を修正基準文章として利用することによって、解析対象文章を、関連性のある意味内容の文章に修正することができる。また、解析対象文章に記述された日時又は作成された日時に近い文章集合を修正基準文章として利用することによって、文章作成時の流行を反映した文章に修正することができる。これにより、異なる意味内容の文章に修正されたり、過剰修正されたりすることなく、文章を修正することができる。

【0082】

また、修正基準文章が、形態素解析の精度が高い一定の文章(例えば新聞記事のみ)である場合と比較して、修正後の文章における未知語の割合を減少させることができ、文章修正の精度を上げることができる。

【0083】

前述した本発明の種々の実施形態について、本発明の技術思想及び見地の範囲の種々の変更、修正及び省略は、当業者によれば容易に行うことができる。前述の説明はあくまで例であって、何ら制約しようとするものではない。本発明は、特許請求の範囲及びその均

10

20

30

40

50

等物として限定するものにのみ制約される。

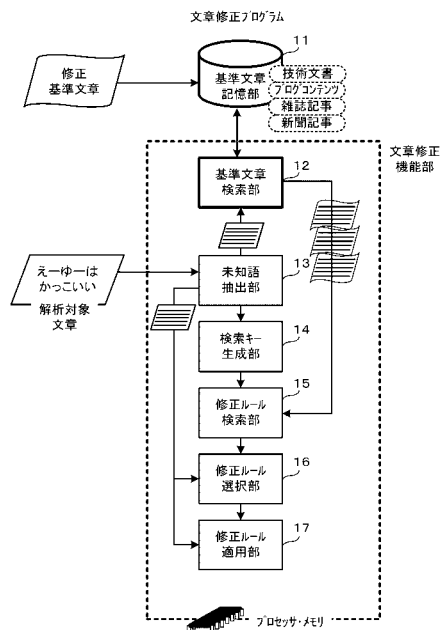
【符号の説明】

【 0 0 8 4 】

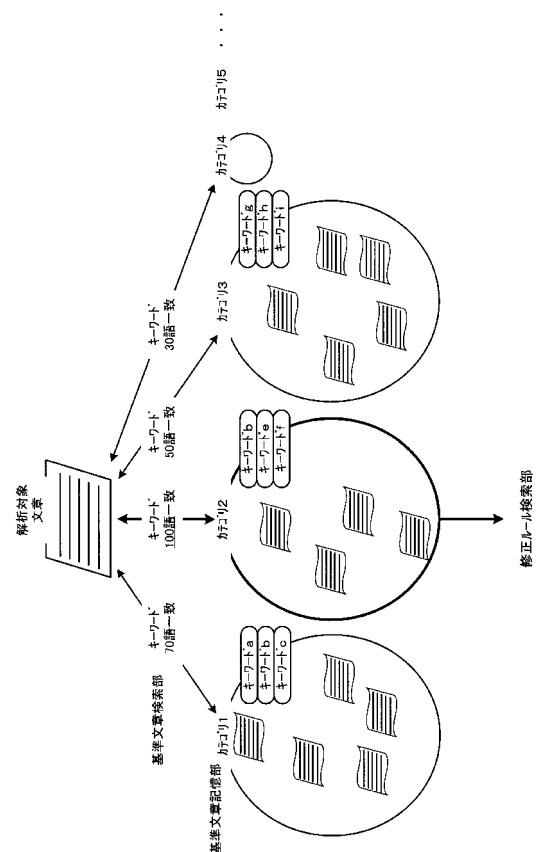
- 1 文章修正プログラム
- 1 1 基準文章記憶部
- 1 2 基準文章検索部
- 1 3 未知語抽出部
- 1 4 検索キー生成部
- 1 5 修正ルール検索部
- 1 6 修正ルール選択部
- 1 7 修正ルール適用部
- 2 文章解析サーバ
- 2 0 通信インタフェース部
- 2 1 基準文章入力部
- 2 2 解析対象文章入力部
- 2 3 文章修正機能部
- 3 W e bサーバ
- 4 投稿用端末

10

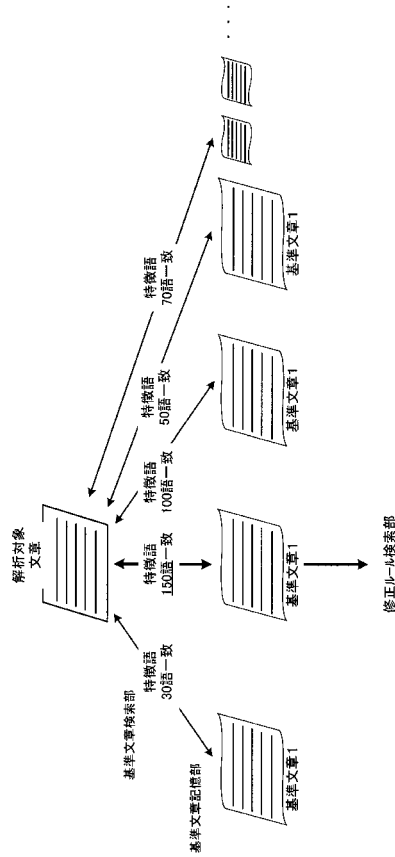
【 図 1 】



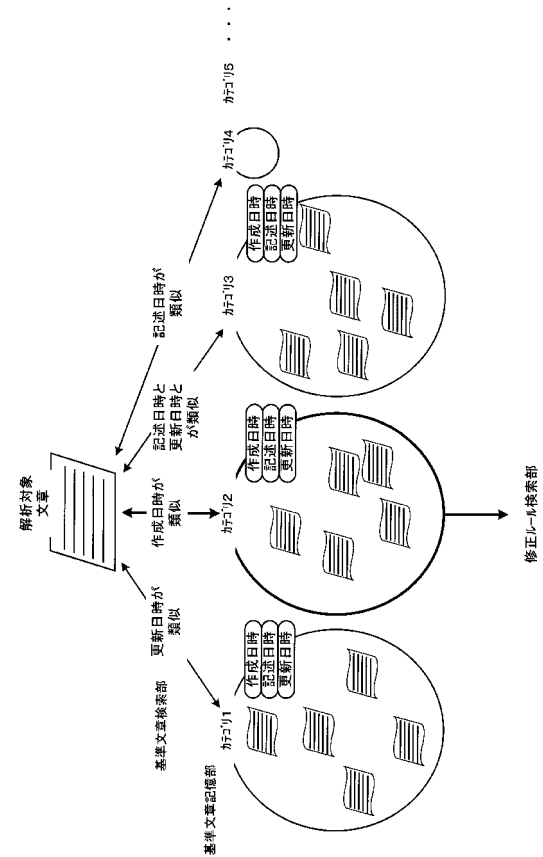
【圖 2】



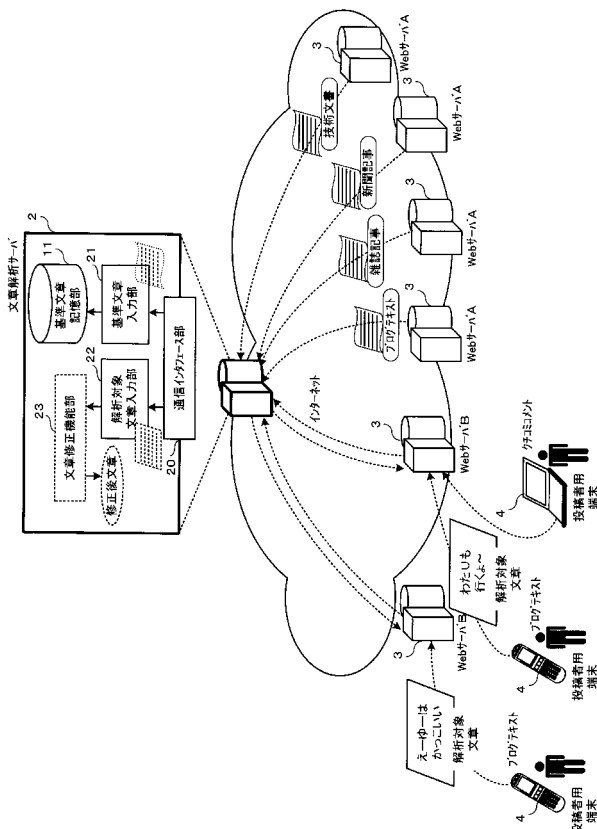
【 図 3 】



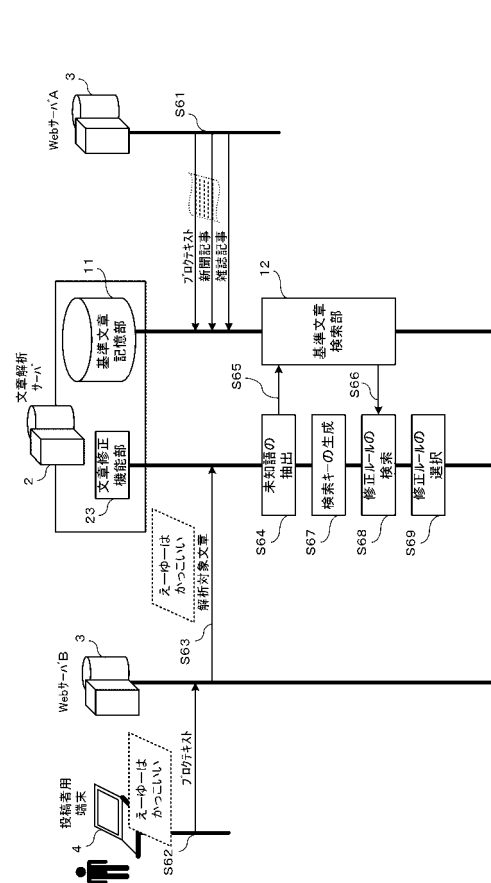
【 図 4 】



【 図 5 】



【 図 6 】



## フロントページの続き

(51)Int.Cl. F I  
G 0 6 F 17/30 3 5 0 C

(72)発明者 滝嶋 康弘  
埼玉県ふじみ野市大原二丁目1番15号 株式会社KDDI研究所内

審査官 長 由紀子

(56)参考文献 特開平09-153034(JP,A)  
特開平05-334293(JP,A)  
大鹿 広憲 外3名, Googleを活用した英作文支援システムの構築, DEWS2005論  
文集 [online], 日本, (社)電子情報通信学会データ工学研究専門委員会, 200  
5年 5月 2日, DEWS2005 4B-i8  
日野 浩平 外5名, ウェブからの関連語収集手法を用いた専門用語の訳語推定, 言語処理学会  
第11回年次大会発表論文集, 日本, 言語処理学会, 2005年 3月15日, p.21-24  
池田 和史 外3名, ブログにおける表記の揺れを修正するためのルール自動作成システムの提  
案, 第71回(平成21年)全国大会講演論文集(2) 人工知能と認知科学, 日本, 社団法人  
情報処理学会, 2009年 3月10日, p.2-79~2-80

(58)調査した分野(Int.Cl., DB名)  
G 0 6 F 1 7 / 2 0 - 2 8  
G 0 6 F 1 7 / 3 0