

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号  
特許第4717821号  
(P4717821)

(45) 発行日 平成23年7月6日(2011.7.6)

(24) 登録日 平成23年4月8日(2011.4.8)

(51) Int.Cl.

F I

G O 6 F 17/30 (2006.01)

G O 6 F 17/28 (2006.01)

G O 6 F 17/30 3 2 O C

G O 6 F 17/30 1 7 O A

G O 6 F 17/28 Z

請求項の数 15 (全 24 頁)

(21) 出願番号	特願2006-533909 (P2006-533909)	(73) 特許権者	502208397
(86) (22) 出願日	平成16年9月13日 (2004. 9. 13)		グーグル インコーポレイテッド
(65) 公表番号	特表2007-507796 (P2007-507796A)		アメリカ合衆国 カリフォルニア州 9 4
(43) 公表日	平成19年3月29日 (2007. 3. 29)		0 4 3 マウンテン ビュー アンフィシ
(86) 国際出願番号	PCT/US2004/029772		アター パークウェイ 1 6 0 0
(87) 国際公開番号	W02005/033967	(74) 代理人	100078282
(87) 国際公開日	平成17年4月14日 (2005. 4. 14)		弁理士 山本 秀策
審査請求日	平成19年8月20日 (2007. 8. 20)	(74) 代理人	100062409
(31) 優先権主張番号	10/676, 724		弁理士 安村 高明
(32) 優先日	平成15年9月30日 (2003. 9. 30)	(74) 代理人	100113413
(33) 優先権主張国	米国 (US)		弁理士 森下 夏樹
		(72) 発明者	ミッタル, ビブ
			アメリカ合衆国 カリフォルニア 9 4 0
			8 7, サニーベール, エルソナ ドラ
			イブ 1 3 2 7
			最終頁に続く

(54) 【発明の名称】 ターゲットページとは異なる文字セットおよび／または言語で書かれたクエリを使用する検索のための方法

(57) 【特許請求の範囲】

【請求項 1】

複数のサーバと1つの入力デバイスとを含むシステムにおける方法であって、該複数のサーバの各々はプロセッサを含み、該方法は、

該複数のサーバのうちの1つのサーバが、第1のフォーマットで書かれ、所定の言葉を含む第1のアンカーテキストのセットを識別することと、

該複数のサーバのうちの該1つのサーバが、該第1のアンカーテキストのセットが指す文書のセットを識別することと、

該複数のサーバのうちの該1つのサーバが、第2のフォーマットで書かれ、該識別された文書のセットを指す第2のアンカーテキストのセットを識別することと、

該複数のサーバのうちの該1つのサーバの該プロセッサが、該第1のフォーマットにおける該所定の言葉の表示が該第2のフォーマットにおける該所定の言葉の表示に対応することを決定するために該第2のアンカーテキストのセットを分析することであって、該第2のアンカーテキストのセットを分析することは、該所定の言葉が該第2のアンカーテキストのセット内の言葉に対応する確率を計算することを含む、こととを包含する、方法。

【請求項 2】

前記第1のフォーマットは第1の文字セットを含み、前記第2のフォーマットは第2の文字セットを含む、請求項1に記載の方法。

【請求項 3】

前記第 1 のフォーマットは第 1 の言語を含み、前記第 2 のフォーマットは第 2 の言語を含む、請求項 1 に記載の方法。

【請求項 4】

前記第 2 のアンカーテキストのセットを分析することは、該第 2 のアンカーテキストのセットで最も頻繁に現われる言葉を識別することと、該第 2 のフォーマットにおいて、該最も頻繁に現われる言葉を前記所定の言葉の表示として指定することを含む、請求項 1 に記載の方法。

【請求項 5】

前記確率は、ベイズ法、ヒストグラムスムージング、カーネルスムージング、および縮小推定量のうちの少なくとも一つを用いて得られる、請求項 1 に記載の方法。

10

【請求項 6】

前記所定の言葉が前記第 2 のアンカーテキストのセットにおける言葉に対応する前記確率は、該第 2 のアンカーテキストのセットにおける該言葉の発生回数を、該第 2 のアンカーテキストのセットにおける全言葉の総発生回数で割ることによって得られる、請求項 1 に記載の方法。

【請求項 7】

前記第 2 のアンカーテキストのセットを分析することは、前記所定の言葉が該第 2 のアンカーテキストのセットにおけるそれぞれの言葉に対応する確率を計算することを包含する、請求項 1 に記載の方法。

【請求項 8】

20

前記第 2 のアンカーテキストのセットを分析することは、該第 2 のアンカーテキストのセットにおいて最も頻繁に現われる言葉を識別することを包含する、請求項 1 に記載の方法。

【請求項 9】

前記第 1 のフォーマットは、ローマ字、r o m a j a、およびピンインから成る群から選択され、前記第 2 の文字セットは、カタカナ、ひらがな、漢字、ハングル、ハンジャ、および伝統的な中国文字から成る群から選択される、請求項 2 に記載の方法。

【請求項 10】

前記文書がウェブページを含む、請求項 1 に記載の方法。

【請求項 11】

30

前記複数のサーバのうちの前記 1 つのサーバが、前記第 1 のフォーマットで書かれ、前記所定の言葉を含むクエリを前記入力デバイスから得ることと、

該複数のサーバのうちの該 1 つのサーバが、少なくとも部分的に前記分析することに基づいて、該クエリを前記第 2 のフォーマットに変換することと、

該複数のサーバのうちの該 1 つのサーバが、データベースで、該変換されたクエリに応じた、該第 2 のフォーマットで書かれた情報を検索することとをさらに包含する、請求項 1 に記載の方法。

【請求項 12】

前記第 1 のアンカーテキストのセットを識別すること、前記文書のセットを識別すること、前記第 2 のアンカーテキストのセットを識別すること、前記分析すること、前記得ること、前記変換すること、および前記検索することは、列挙された順序で実行される、請求項 11 に記載の方法。

40

【請求項 13】

コンピュータプログラムを格納したコンピュータ読み取り可能格納媒体であって、該コンピュータプログラムは、プロセッサを含むコンピュータシステムによって実行される場合に該コンピュータシステムに動作を実行させるように操作可能な命令を含み、該動作が、

第 1 のフォーマットで書かれ、所定の言葉を含む第 1 のアンカーテキストのセットを識別することと、

該第 1 のアンカーテキストのセットが指すウェブページのセットを識別することと、

50

第2のフォーマットで書かれ、該識別されたウェブページのセットを指す第2のアンカーテキストのセットを識別することと、

該プロセッサが、該第1のフォーマットにおける該所定の言葉の表示が該第2のフォーマットにおける該所定の言葉の表示に対応する確率を計算することとを包含する、コンピュータ読み取り可能格納媒体。

【請求項14】

前記コンピュータシステムによって実行される場合に該コンピュータシステムに動作を実行させるように操作可能な命令をさらに含み、該動作が、

前記プロセッサが、少なくとも部分的に、前記第2のフォーマットにおける前記所定の言葉の表示を含む検索結果をユーザーが選択することの、該第2のフォーマットにおける該所定の言葉の他の表示を含む検索結果をユーザーが選択することに対する比率に基づいて、前記第1のフォーマットにおける該所定の言葉の表示が該第2のフォーマットにおける該所定の言葉の表示に対応する確率を修正することを包含する、請求項13に記載のコンピュータ読み取り可能格納媒体。

10

【請求項15】

前記確率は、少なくとも部分的に、ベイズ法、ヒストグラムスムージング、カーネルスムージング、および縮小推定量のうちの少なくとも一つを用いて計算される、請求項13に記載のコンピュータ読み取り可能格納媒体。

【発明の詳細な説明】

【技術分野】

20

【0001】

(関連出願の参照)

本出願は、2000年12月26日に出願され、「METHODS AND APPARATUS FOR PROVIDING SEARCH RESULTS IN RESPONSE TO AN AMBIGUOUS SEARCH QUERY」と題された、米国特許出願シリアル番号第09/748,431号の一部継続であり、2000年7月6日に出願され、「DATA ENTRY AND SEARCH FOR HANDHELD DEVICES」と題された、米国特許仮出願シリアル番号第60/216,530号の優先権を、米国特許法第119条(e)に基づき主張し、その両方は、それらの全体においてここに援用される。

30

【0002】

(発明の分野)

本発明は一般に、情報検索に関する。より詳細には、検索される文章の少なくとも一部の文字セットまたは言語とは異なる文字セットまたは言語において書かれたクエリを使用して検索を実行するためのシステムおよび方法が開示される。

【背景技術】

【0003】

多くの検索エンジンは、エンドユーザが従来のキーボードなどのようなものを用いて検索クエリを入力するという想定の下で動作し、そこで、英数字の入力は難しいことではない。小さなデバイスがより一般的になってはいるが、しかしながら、この想定はいつも有効とは限らない。例えば、ユーザは、WAP(ワイヤレス・アプリケーション・プロトコル)規格をサポートする携帯電話を使用して、検索エンジンにクエリし得る。携帯電話などのデバイスは通常、データ入力インターフェースを有し、ユーザによる特定のアクション(例えばキーを押すなど)が一つ以上の英数字文字に対応し得る。WAP構成の詳細は、<http://www1.wapforum.org/tech/documents/SPEC-WAPArch-19980439.pdf>(「WAP 100 Wireless Application Protocol Architecture Specification」)にて利用可能である。

40

【0004】

通常の場合、WAPユーザは、検索クエリのページにナビゲートされ、ユーザが検索ク

50

エリを入力するフォームを提示される。従来の方法では、ユーザは、多数のキーを押して、特定の文字を選択し得る。標準の電話のキーパッドでは、ユーザは、例えば、文字「b」を選択する場合、「2」のキーを2回押す。または、文字「s」を選択したい場合は、「7」のキーを4回押す。従って、「ben smith」というクエリを入力するには、ユーザは通常、2 2 3 3 6 6 0 7 7 7 7 6 4 4 4 8 4 4 という一連のキーを押して入力する必要があり、以下のような文字に対応する。

2 2    b  
3 3    e  
6 6    n  
0    スペース  
7 7 7 7    s  
6    m  
4 4 4    i  
8    t  
4 4    h

10

ユーザが検索リクエストを入力した後、検索エンジンは、ユーザから文字を受け取り、あたかも、ユーザが従来のキーボードを用いて、デスクトップのブラウザからリクエストを受け取ったかのように、同様の方法で処理する。

【0005】

前述の例から理解できるように、データ入力のこの形式は、「ben smith」に対応する9つの英数字文字（スペースを含む）を入力するために18回もキー入力が必要とする点で、非効率的である。

20

【0006】

同様の困難さは、ターゲットでない(non-target)言語のキーボードを用いてクエリをタイピングする場合に生じ得る。例えば、日本語のテキストは、ひらがな、カタカナ、および漢字などを含む様々な異なる文字セットを用いて表現され得、そのどれもが、ローマ字(Roman alphabet)に基づいた通常のASCIIキーボードを用いて容易に入力されるものではない。そのような状況において、ユーザはしばしば、日本の徳島市所在のJust System Corp.によって製造されたIchitarron(登録商標)などのようなワードプロセッサのソフトを使用し、romaji(日本語における音声的なローマ字(Roman alphabet)の表現)で書かれたテキストを、カタカナ、ひらがな、および漢字に変換することができる。ワードプロセッサソフトを使用し、ユーザはローマ字でクエリをタイピングし、次いで、ワードプロセッサのスクリーンから変換されたテキストを、ブラウザの検索ボックスへとカットアンドペーストする。このアプローチの不利な点は、相対的に遅く、面倒であり得、ユーザがワードプロセッサのコピーにアクセスすることが要求されるゆえ、コストの制約やメモリの制約などのために、ふさわしいとはいえない。

30

【0007】

それゆえ、曖昧な検索クエリに応じて、適切な検索結果を提供する方法および装置が必要とされるのである。

40

【発明の開示】

【課題を解決するための手段】

【0008】

具体化され、ここで広く記載される本発明と合致する方法および装置は、曖昧な検索クエリに応じた適切な検索結果を提供する。本発明と合致し、そのような方法は、ユーザからの一連の曖昧な情報構成要素を受け取ることを含む。その方法は、曖昧な情報構成要素を、より曖昧でない情報にマッピングする、マッピング情報を含む。このマッピング情報は、一連の曖昧な情報構成要素を、一つ以上の対応する一連の、より曖昧でない情報構成要素に変換するために使用される。一つ以上のこれらの一連の、より曖昧でない情報構成要素は検索エンジンへの入力として提供される。その検索結果は検索エンジンから得られ

50

、ユーザに提示される。

【 0 0 0 9 】

付け加えて、システムおよび方法は、検索される文書の少なくとも一部の文字セットまたは言語とは異なる文字セットにて表されたクエリを用いて検索を実行することが開示される。本発明の実施形態により、ユーザは、標準の入力デバイス（例えば、A S C I I キーボード）を用いてクエリをタイプすることができ、クエリをサーバにおいて適切な形式に変換させることができ（たとえば、ローマ字で書かれたクエリをカタカナ、ひらがな、および／または漢字に変換する）、ならびに、変換された形式に基づいて、検索結果を受け取ることができる。

【 0 0 1 0 】

本発明は、プロセス、装置、システム、デバイス、方法、または、コンピュータ可読格納媒体、搬送波、またはコンピュータネットワークなどのコンピュータ可読媒体を含み、多様な方法においてインプリメントされ得ることは理解されるべきであり、プログラムの命令は、光学式または電気の通信線を介して送信される。いくつかの発明の実施形態は以下に記載される。

【 0 0 1 1 】

一実施形態において、方法は、クエリの言葉を、一つの言語および／または文字セットから別のものへと、自動的に変換することが記載される。所定のクエリの言葉を含むアンカーテキストの第1のセットが識別され、それはアンカーテキストが提示する文書（例えばウェブページなど）のセットである。次いで、第2のフォーマットで書かれ、同じ文書のセットを提示するアンカーテキストの第2のセットが識別される。アンカーテキストの第2のセットは、次いで、分析され、第1のフォーマットにおける所定のクエリの言葉の表示が、第2のフォーマットにおける所定のクエリの言葉の表示に対応する確率を得る。

【 0 0 1 2 】

別の実施形態において、確率辞書が作成され、第1のフォーマット（例えば、言語および／または文字セット）で書かれた言葉を、第2のフォーマット（例えば、別の言語および／または文字セット）にマッピングする。確率辞書は、第1のフォーマットで書かれたクエリを第2のフォーマットに変換するために使用される。変換されたクエリは、次いで、検索を実行するために使用され、その結果は、ユーザに戻される。一部の実施形態において、検索結果を用いたユーザの相互作用は、監視され得、確率辞書における確率を更新するために使用される。また、一部の実施形態において、クエリ自体は、検索に先立って、代替的な言語および／または文字セットのマッピングを含むように拡張され得る。

【 0 0 1 3 】

さらなる別の実施形態において、確率辞書を作成する方法が記載される。確率辞書は、第1のフォーマットにおける言葉を第2のフォーマットに変換するために使用され得る。辞書は、アンカーテキストまたはその言葉を含む他のデータを識別することによって、好ましくは言葉毎に作成される。次に、アンカーテキストまたは他のデータに連係される（*a l l i g n e d w i t h*）データは分析され、第1のフォーマットにおける所定の言葉が、第2のフォーマットにおける一つ以上の言葉にマッピングされる確率を決定する。

【 0 0 1 4 】

さらなる別の実施形態において、第1の言語または文字セットに提供されたクエリは、一つ以上のクエリの言葉を含み、第1の言語または文字セットで書かれたアンカーテキストと、第1のアンカーテキストに対応し、第2の言語または文字セットで書かれたアンカーテキストとを比較することによって、第2の言語または文字セットに変換される。

【 0 0 1 5 】

別の実施形態において、コンピュータプログラム製品は、第1のフォーマットで書かれた言葉を第2のフォーマットに変換するために提供される。コンピュータプログラム製品は、コンピュータシステムに、連係されたアンカーテキストを識別させ、第1のフォーマットにおける所定の言葉の表示が、第2のフォーマットにおける一つ以上の言葉に対応する確率を決定させるように動作可能である。

10

20

30

40

50

## 【 0 0 1 6 】

別の実施形態において、方法は、曖昧なクエリを用いて検索を実行するために提供される。ユーザが第1のフォーマットにおいてクエリを入力する場合、それは、第2のフォーマットで書かれた一つ以上の変形の一群に変換される。次いで、検索は、変換された変形を用いて実行され、応答の情報は、ユーザに戻される。例えば、第1のフォーマットは、電話キーパッドを用いて入力された一連の数を含み得、第2のフォーマットは、英数字のテキスト（例えば、英語、ローマ字、`roma ja`、ピンインなど）を含み得る。一部の実施形態において、一つ以上の変形の群は、所定の語彙に現れない、および/または、所定の低い確率の文字の組み合わせを含む、変換された変形を除去することによって選択される。一部の実施形態において、確率辞書は、検索が実行される前に、一つ以上の変形の群を、第3のフォーマットに変換する。例えば、確率辞書は、ローマ字、`roma ja`、またはピンインの一つ以上の変形の群を、漢字、カタカナ、ひらがな、ハングル、ハンジャ、または伝統的な中国文字（`traditional Chinese character`）に変換するために使用され得、検索は、次いで、変換された変形を用いて実行され得る。

10

## 【 0 0 1 7 】

本発明のこれらおよび他の特徴および利点は、以下の詳細な記載、ならびに、本発明の原理の例によって例示された、添付された図面に、さらに詳細に提示されている。

## 【 発明を実施するための最良の形態 】

## 【 0 0 1 8 】

20

添付された図面は、この明細書にて援用され、その一部として構成され、本発明の実施形態を例示し、記載とともに、本発明の利点および原理を説明するのに役立つ。

## 【 0 0 1 9 】

添付された図面にて例示される本発明の実施形態を詳細に参照する。同様の数字は、図面や以下に続く記載を通して、同様の部分を示す。以下に続く記載は、当業者が本発明を利用することができるように提示される。特定の実施形態および応用の記載は例としてのみ提供されるのであり、様々な修正は当業者にとって容易に明白である。例えば、多くの例がインターネットのウェブページに記載されているが、本発明の実施形態は、本、新聞、雑誌などの文章および/または情報の他のタイプを検索するために使用され得る。同様に、例示のために、日本語のテキストをローマ字からカタカナ、ひらがなおよび/または漢字へと変換されることが記載されるが、当業者に明らかなように、本発明のシステムおよび方法は、任意の適切な変換へと応用され得る。例えば、限定なしに、本発明の実施形態は、一部の他のフォーマット（例えば、ピンインやローマ字など）において受け取られるクエリに基づき、伝統的な漢字、または韓国のハングル文字またはハンジャ文字にて書かれたテキストを検索するために用いられ得る。ここで記載される一般的な原理は、本発明の趣旨および範囲から逸脱することなく他の実施形態および応用に適用され得る。したがって、本発明は、ここで開示される原理および特徴に合致する多数の代替、修正、および均等物を含み、最も広い範囲に従うものである。明瞭さのために、本発明に関連する領域で既知である技術的事項に関する詳細は、本発明を不必要に曖昧にしないように、詳細に記載されていない。

30

40

## 【 0 0 2 0 】

## A . 概説

本発明に合致する方法および装置により、ユーザは、曖昧な検索クエリを提出し、場合によっては明確にされた検索結果を受け取ることが可能である。一実施形態において、標準の電話のキーパッドのユーザから受け取る一連の数が、場合によってはそれに対応する英数字のシーケンスのセットに変換される。これらの対応する英数字のシーケンスは、ブール式の「OR」結果を使用し、従来の検索エンジンへ入力として提供される。この方法において、検索エンジンは、ユーザが興味を持ちそうなものに対する検索結果を制限するのに役立つ。

## 【 0 0 2 1 】

50

## B．構成

本発明に合致する方法および装置がインプリメントされ得る、システム１００が図１に示される。システム１００は、ネットワーク１４０を介して、多数のサーバ１２０および１３０に接続される多数のクライアントデバイス１１０を含み得る。ネットワーク１４０は、ローカルエリアネットワーク（ＬＡＮ）、ワイドエリアネットワーク（ＷＡＮ）、公衆交換電話網（ＰＳＴＮ）などのような電話網、インターネット、またはネットワークの組み合わせを含み得る。２つのクライアントデバイス１１０および３つのサーバ１２０および１３０は、単純に、ネットワーク１４０に接続されるように例示される。実際には、ほぼ同数のクライアントデバイスおよびサーバが存在し得る。たま、一部の案件においては、クライアントデバイスは、サーバ機能を実行し得、サーバはクライアントデバイス機能を実行し得る。

10

### 【００２２】

クライアントデバイス１１０は、メインフレーム、ミニコンピュータ、パーソナルコンピュータ、ラップトップ、ＰＤＡ（携帯情報端末）などのようなデバイスを含み得、ネットワーク１４０に接続可能である。クライアントデバイス１１０は、ネットワーク１４０を介してデータを送信し、あるいは、有線、無線、または光学式の接続を介してネットワーク１４０からデータを受信する。

### 【００２３】

図２は、本発明と合致する例示的なクライアントデバイス１１０を図示する。クライアントデバイス１１０は、バス２１０、プロセッサ２２０、メインメモリ２３０、読み出し専用メモリ（ＲＯＭ）２４０、記憶装置２５０、入力デバイス２６０、出力デバイス２７０、および通信インターフェース２８０を含み得る。

20

### 【００２４】

バス２１０は、一つ以上の従来のバスを含み得、クライアントデバイス１１０の間で通信を可能にする。プロセッサ２２０は、従来のタイプのプロセッサまたはミニコンピュータを含み得、命令を解釈し実行する。メインメモリ２３０はランダムアクセスメモリ（ＲＡＭ）または他のタイプのダイナミック記憶装置を含み得、プロセッサ２２０による実行のための情報および命令を格納する。ＲＯＭ２４０は、従来のＲＯＭデバイスまたは他のタイプのスタティック記憶装置を含み得、プロセッサ２２０が使用するためのスタティックな情報および命令を格納する。記憶装置２５０は磁気および／または光学式の記憶媒体、ならびにそれに対応するドライブを含み得る。

30

### 【００２５】

入力デバイス２６０は、キーボード、マウス、ペン、音声認識および／または生体認識メカニズムなどのような一つ以上の従来のメカニズムを含み得、それによって、ユーザは、クライアントデバイス１１０へ情報を入力することが可能である。出力デバイス２７０は、一つ以上の従来のメカニズム（ディスプレイ、プリンタ、スピーカなど）を含み得、ユーザに情報を出力する。通信インターフェース２８０は、任意の送受信器のようなメカニズムを含み得、クライアントデバイス１１０が他のデバイスおよび／またはシステムと通信することが可能となる。例えば、通信インターフェース２８０は、ネットワーク１４０などのようなネットワークを介して、別のデバイスまたはシステムと通信するためのメカニズムを含み得る。

40

### 【００２６】

以下で詳細に記載するように、本発明と合致するクライアントデバイス１１０は、所定の検索に関連する動作を実行する。クライアントデバイス１１０は、メモリ２３０などのコンピュータ可読媒体に含まれるソフトウェアの命令を実行するプロセッサ２２０に回答した動作を実行し得る。コンピュータ可読媒体は一つ以上のメモリデバイスおよび／または搬送波として定義され得る。ソフトウェアの命令は、データ記憶装置２５０などのような別のコンピュータ可読媒体から、または、通信インターフェース２８０を介して別のデバイスから、メモリ２３０へと読み出され得る。メモリ２３０に含まれるソフトウェアの命令により、プロセッサ２２０は、以下で記載される、検索に関連する動きを実行する。

50

あるいは、ハードウェアに組み込まれている回路は、ソフトウェアの命令の代わりに、またはソフトウェアの命令と組み合わせられて使用され得、本発明と合致する処理をインプリメントし得る。したがって、本発明は、特定のハードウェアに組み込まれている回路とソフトウェアとの任意の組み合わせに限定されるわけではない。

【0027】

サーバ120および130は、メインフレーム、ミニコンピュータ、またはパーソナルコンピュータなどの、一つ以上のコンピュータシステムのタイプを含み得、ネットワーク140と接続することができ、サーバ120および130は、クライアントデバイス110と通信することができる。代替的な実施において、サーバ120および130は、一つ以上のクライアントデバイス110と直接に接続するメカニズムを含み得る。サーバ120および130は、ネットワーク140を介してデータを送信し得、あるいは、有線、無線、または光学式の接続を介して、ネットワーク140からデータを受信し得る。

10

【0028】

サーバは、クライアントデバイス110に対して、図2を参照し上記されたのと同様な方法にて、構成され得る。本発明と合致する実施において、サーバ120は、クライアントデバイス110によって使用可能である検索エンジン125を含み得る。サーバ130は、クライアントデバイス110によってアクセス可能である文書（またはウェブページ）を格納し得る。

【0029】

C. 構成動作

20

図3は、3つの文書を表す図を示し、サーバ130のうちの一つにおける例示として格納され得る。

【0030】

第1の文書（文書1）は、「car repair」および「car rental」の2つのデータ入力を含み、底に「3」という数がある。第2の文書（文書2）は「video rental」というデータ入力を含む。第3の文書（文書3）は、「wine」、「champagne」、および「bar items」という3つのデータ入力を含み、ならびに、文書2へのリンク（または参照）を含む。

【0031】

例示を単純にするために、図3に示される文書は、英数字文字列の情報（例えば、「car」、「repair」、「wine」など）のみを含む。しかしながら、当業者は、他の状況において、文書は、音声的、視聴覚的な情報などといった、他のタイプの情報を含み得ることを理解する。

30

【0032】

図4aは、図3に示された文書に基づき、従来の英数字のインデックスを示す。インデックスの第1の列は英数字のリストを含み、第2の列は、それらの言葉に対応する文書のリストを含む。英数字の「3」などの一部の言葉は、一つの文書（この場合は文書1）に対応する（にある）。「rental」などの他の言葉は、多数の文書（この場合は文書1および2）に対応する。

【0033】

40

図4bは、検索エンジン125のような従来の検索エンジンが、図4aにて示されたインデックスをどのように使用し、英数字の検索クエリに応じた検索結果を提供するのかを示す。英数字のクエリは、任意の従来の技術を用いて生成され得る。図示のために、図4bは2つの英数字のクエリ、「car」および「wine」を表す。従来のアプローチの下では、検索エンジン125は、「car」などの英数字のクエリを受け取り（ステージ410）、その英数字のインデックスを用いて、どの文書がそのクエリに対応するのかを決定する（ステージ420）。この例において、従来の検索エンジン125は、図4aにて示されたインデックスを用い、「car」が文書1に対応し、検索結果としてユーザに、文書1（またはそれへの参照）を戻す。同様に、従来の検索エンジンは、「wine」が文書3に対応し、ユーザに文書3（またはそれへの参照）を戻すことを決定する（ステ

50



ージ 4 3 0 )。

【 0 0 3 4 】

図 5 a は、本発明と合致し、個々に図 3 および図 4 a に示された文書およびインデックスに基づき、数字の検索クエリに応じた検索結果を提供する好ましい技術のフロー図を示す。理解し易い例示のために、図 5 a は標準の電話端末のマッピングに基づき、数字のクエリを処理する特定の技術を記載する。しかし、当業者は、本発明に合致する他の技術もまた使用され得ることを理解する。

【 0 0 3 5 】

ステージ 5 1 0 において、シーケンス「 2 2 7 」( 数字の構成要素「 2 」、「 2 」、および「 7 」からなる ) がユーザから受け取られる。ステージ 5 2 0 において、数字の構成要素が文字へとマッピングされる方法についての情報が得られる。ユーザが標準の電話キーボードから情報を入力したと想定し、このマッピング情報が図 5 b に示される。図 5 b に示されるように、文字「 a 」、「 b 」、および「 c 」はそれぞれ、「 1 」にマッピングされ、文字「 p 」、「 q 」、および「 r 」はそれぞれ、「 7 」にマッピングされている。

【 0 0 3 6 】

ステージ 5 3 0 において、このマッピング情報を用い、シーケンス「 2 2 7 」は、その英数字の相当物に変換される。図 5 b に示される情報に基づいて、シーケンス「 2 2 7 」に対応する文字の可能な組み合わせは 3 6 通りあり、その組み合わせは、 a a p 、 b a p 、 c a p 、 a b p 、 b b p 、 . . . b a r . . . c a r . . . c c s などを含む。数字が可能な組み合わせ ( 例えば「 a a 7 」 ) に含まれる場合、 8 0 通りの可能な組み合わせが存在する。全ての可能な英数字の相当物を生成するよりもむしろ、一部の語彙に基づき、生成された相当物を限定するのが望ましい。例えば、辞書、または以前の検索クエリの検索エンジンログなどに存在する英数字の相当物のみを生成することが望ましい。あるいは、既知の統計的な技術 ( 例えば、所定の言葉と一緒に現れる確率など ) を用いることによって英数字の相当物を限定することが望ましい。

【 0 0 3 7 】

ステージ 5 4 0 において、これらの英数字の相当物は、論理「 O R 」動作を用いて、図 4 a および図 4 b を参照して記載されたような、従来の検索エンジンへの入力データとして提供される。例えば、検索エンジンへ提供される検索クエリは、「 a a p O R b a p O R c a p O R a b p . . . O R b a r . . . O R c a r 」であり得る。全ての可能な英数字の相当物が検索クエリに提供され得るが、サブセットは、その代わりに、意図されない相当物を除去する従来の技術を用いて使用され得る。例えば、ユーザは、言葉の使用についての確率的な情報を引き出す技術を用いて、可能な組み合わせのより狭いリストを生成することが可能である。すなわち、ユーザは、「 q u 」で始まる組み合わせを含む ( 好む ) が「 q t 」で始まる組み合わせを無視することが可能である。

【 0 0 3 8 】

ステージ 5 5 0 において、検索結果は検索エンジンから得られる。なぜなら、「 a a p 」や「 a b p 」などの言葉は検索エンジンのインデックスには存在せず、それらは効果的に無視されるからである。実際には、図 4 b に示されたインデックス内に含まれた言葉は「 c a r 」および「 b a r 」のみであり、戻ってきた検索結果は、文書 1 および文書 3 を参照するもののみである。ステージ 5 6 0 において、これらの検索結果はユーザに提示される。その検索結果は、検索エンジンによって提供されたのと同じ順序にて提示され得るか、または、ユーザの言語などの検討材料に基づき、記録され得る。ユーザが「 b a r 」という言葉を含む文書のみに興味を持っていると想定すると、ユーザは、望んでいない結果 ( 文書 3 ) 、および望んだ結果 ( 文書 1 ) を受け取る。ユーザの利点として、検索クエリを編成する 3 つのキーを押すことが必要とされるのみであるが、これは許容可能な対価であり得る。

【 0 0 3 9 】

図 6 は、本発明に合致し、個々に図 3 および図 4 a に示された文書およびインデックスに基づき、数字の検索クエリに対応する検索結果を提供する好ましい技術の、別のフロー

10

20

30

40

50

図を示す。このフロー図は、受け取られたシーケンスのサイズの増加が、ユーザによって望まれたものへの検索結果をどのように限定するのを助けるのかを例証する。理解し易い例示のために、図6は、標準の電話キーパッドのマッピングに基づいて、数字のクエリを処理する特定の技術を再び記載するが、当業者は、本発明に合致する他の技術が利用され得ることを理解する。

#### 【0040】

ステージ610において、シーケンス「227 48367」（数字の構成要素、「2」、「2」、「7」、「4」、「8」、「3」、「6」、「7」からなる）が、ユーザから受け取られる。説明のために、シーケンス「227」を「数字ワード」（number word）と呼び、全体のシーケンス「227 48367」を「数字フレーズ」（number phrase）と呼ぶ。数字ワードの可能な英数字の相当物を「文字ワード」（letter word）と呼び、数字フレーズの可能な英数字の相当物を「文字フレーズ」（letter phrase）と呼ぶ。

10

#### 【0041】

ステージ620において、数字の構成要素が文字にマッピングされる方法についての情報が得られる。ステージ630において、同様のマッピング情報が、図5bに示されるように使用されると想定し、数字フレーズ「227 48367」は、それに対応する文字フレーズに変換される。図5bに示される情報に基づき、シーケンス「227 48367」に対応する、11664通りの可能な文字フレーズが存在する。

#### 【0042】

20

ステージ640において、これらの文字フレーズは、論理「OR」動作を用いて、図4aおよび図4bを参照して記載された、従来の検索エンジンへの入力データとして提供される。例えば、検索エンジンに提供された検索クエリは、「aap gtdmp 'OR' aap ht dmp '...OR' bar items '...OR' car items」であり得る。全ての可能な文字フレーズが検索エンジンに提供され得るが、サブセットは、その代わり、意図されていない文字フレーズを除去するために従来の技術を用いて、使用され得る。

#### 【0043】

ステージ650において、検索結果は検索エンジンから得られる。なぜなら、多くの検索エンジンが、ソートされた正確なフレーズを含むそれらの文書を上位にランクさせるように設計されており、文書3は、最上位にランクされた検索結果であるからである（つまり、正確なフレーズ、「bar items」を含むからである）。例における文書で、ステージ620にて生成された他の文字フレーズのうちの一つを含む文書はない。さらに、多くの検索結果は、フレーズの個々の部分を含む検索結果を減らす（除去する）が、全体のフレーズではない。例えば、文書1は、それが、「car」という文字ワードを含むゆえに減らされ（除去され）、その文字ワードは文字フレーズの第1の部分に対応するが、それは、文字フレーズの第2の部分に対応する任意の文字ワードを含まない。最後に、「aap ht dmp」などの文字フレーズは効果的に無視される。というのは、それらは、検索エンジンのインデックスに存在する文字ワードを含まないからである。

30

#### 【0044】

40

ステージ660において、検索結果がユーザに提示される。例において、ユーザに示された第1の結果は文書3であり、それは、ユーザのクエリに最も適切なものである。文書1は、可能な文字フレーズのうちの一つを含まないゆえに、共に除去され得る。この方法において、ユーザは、最も適切な検索結果が提供される。

#### 【0045】

図5および図6を参照した上記は、数字の情報を受け取り、それを英数字の情報にマッピングすることを参照してなされるが、当業者は、他の実施が本発明と合致して可能であることを理解する。例えば、ユーザによって押されたキーに対応する数のシーケンスを受け取る代わりに、受け取られたシーケンスは、ユーザによって押されたキーに対応する第1の文字からなり得る。つまり、「227」を受け取る代わりに、受け取られたシーケン

50

スは「a a p」であり得る。本発明と合致し、ステージ530または630において生成された、その相当する文字シーケンスは、「a a p」に対応する他の文字シーケンス（例えば「b a r」）であり得る。実際には、受け取られたシーケンスは、音声的、視聴覚的、または他の任意の情報構成要素のタイプを含み得る。

【0046】

シーケンスが受け取られるフォームに関係なく、受け取られたシーケンスは、情報が検索エンジンのインデックスに格納されるフォーマットに対応するシーケンスに変換されるのが、通常は好ましい。例えば、検索エンジンのインデックスが英数字のフォーマットにて記憶される場合、受け取られたシーケンスは英数字のシーケンスに変換されるべきである。

10

【0047】

さらに、情報構成要素の受け取られたシーケンスを変換するために使用されるマッピング技術は、ユーザの入力をデバイスによって生成される情報にマッピングするためのユーザのデバイスにて用いられるのと同じ技術であることが通常は望ましい。しかしながら、ユーザの入力に使用されるのとは異なるマッピング技術が使用されるのが好ましい場合もあり得る。

【0048】

また、本発明の実施形態により、ユーザは、ターゲットではない言語のキーボードを用いて入力された検索を実行することを可能にする。例えば、日本語のテキストを含むウェブページは漢字で書かれる一方、そのページを検索しようとするユーザは、ローマ字のアルファベットに基づいて標準のASCIIキーボード（または携帯電話機）にアクセスを有するのみである。

20

【0049】

図7は、そのような検索を実行する方法を例示する。図7において示されるように、ユーザは、標準の入力デバイス（例えば、ASCIIキーボード、携帯電話機、など）を使用してクエリをタイプし、そのクエリを検索エンジンに送る。そのクエリは、それに対応する文書の一部が書かれている（例えば漢字）文字セットとは異なる文字セット（例えばローマ字）で書かれ得る。検索エンジンは、クエリを受信し（ブロック702）、それを適切な形式に変換し（ブロック704）、従来の検索技術などを使用して、変換されたクエリに対応して、文書の検索を実行する（ブロック706）。次いで検索エンジンは、対応する文書のリスト（および/または、文書のコピー）を、ユーザに戻す（ブロック708）。例えば、図6と関連して上記されたものと同様の方法にて、結果はユーザに戻され得る。

30

【0050】

図7に示されるように、ユーザのクエリは、好ましくは、クライアントとは逆の、検索エンジンサーバ側にて変換され、変換を実行するための、特殊な目的のためのソフトウェアを得る必要から、ユーザを解放させる。しかしながら、他の実施形態において、変換の一部または全ては、クライアント側にて実行され得ることは理解される。付け加えて、一部の実施形態において、クエリは、電話機のキーパッドなどのようなデバイスを使用して入力され得る。そのような実施形態において、初期の数字のクエリは、図5および図6に関連する上記されたマッピング技術を用いて、英数字の形式（例えばローマ字）に変換され得、例えば、低い確率のマッピング（例えば、ローマ字においては生じない文字の組み合わせなどを含むマッピング）を除去するための、語彙および/または確率的な技術の応用を含む。いったん、クエリの英数字の変換が得られると、図7に示されるステップの残りが実行され得る（つまり、704、706、および708）。

40

【0051】

一つの文字セットまたは言語から別のものへの変換（つまり、図7におけるブロック704）は、様々な方法にて実行され得る。一つの技術は、クエリにおけるそれぞれの言葉を、ターゲットの言語または文字セットにおいて対応する言葉にマッピングするために、言葉の意味または変換の従来のスタティックの辞書を使用する。しかしながら、このアプ

50

ローチの問題は、しばしば不正確な結果を生じることである。というのは、言葉はしばしば曖昧であり、クエリはしばしば、短すぎて、この曖昧さを解消する十分な手掛かりを提供しない。例えば、「bank」という言葉は、川の土手、金融機関、飛行機による演習、などを意味し得、理論上、正確に変換することは困難である。付け加えて、辞書が相対的に、小さくなく、および/または頻繁に更新されていない場合、滅多に使用しない言葉、スラング、イディオム、適切な名前などの、検索エンジンが出くわし得る全ての言葉の入力が含まれ得ない。

#### 【0052】

本発明の実施形態はまた、一つの言語または文字セット（例えばASCII）から別のもの（例えば漢字）へ、クエリの変換するために、確率辞書（probabilistic dictionary）を使用して、一部または全ての問題を克服または改善するために使用され得る。好ましい実施形態において、確率辞書は、一つの言葉のセットを別の言葉のセットへとマッピングし、確率をそのマッピングそれぞれに関連付ける。便宜上、「言葉（term）」または「トークン」は、言葉（words）、フレーズ、および/または、スペースを含み得る一つ以上の文字のシーケンスを参照する。

#### 【0053】

図8は、上記されたような確率辞書800の例を示す。図8に示された例示的な確率辞書800は、romaji（日本語のローマ字のアルファベットの表示）で書かれた言葉を、漢字（ローマ字ではない、表意文字ベースの日本語の文字セット）で書かれた言葉にマッピングする。説明を容易にするために、図8は、ローマ字の言葉を、「<term>romaji」とし、漢字の言葉を「<term>kanji」とする。漢字辞典に対する実際のローマ字において、実際のローマ字および漢字の言葉は、図8にて示される英語変換よりも、使用されることが理解される。したがって、図8は、本発明の実施形態の説明を容易にするために提供されるのであり、日本語のテキストの実際の文字および意味を例示しているのではない。

#### 【0054】

辞書800は、様々なローマ字の言葉802のための808、810、812、814のデータ入力を含む。辞書はまた、漢字804におけるそれぞれの言葉の表示を含み、それぞれの表示が正しい場合の、対応する確率806に沿っている。例えば、ローマ字の言葉「bank」は、「steep slope」を意味する言葉に、0.3の確率でマッピングされ得、「financial institution」を意味する言葉には、0.4の確率でマッピングされ得、「airplane maneuver」には0.2の確率でマッピングされ得る。0.1の確率では、その言葉は「その他」にマッピングされ得、それは単に、辞書に存在し得ない言葉に、それぞれの言葉をマッピングさせる、包括的な方法である。

#### 【0055】

再び、図8に示された例は、第1の文字セットまたは言語における所定の言葉（例えば、「bank」など）が、別の文字セットまたは言語において2つ以上の言葉にマッピングされ得ることを例示するように構成されることが理解される。しかしながら、当業者が理解するように、明瞭さのために、図8における特定の例は、英語の言葉や意味を使用する原理を例示しており、「bank」などの言葉の実際のローマ字表示は、例えば、その英語の相当物と同じような方法にて曖昧ではあり得ない（例えば、financial institutionとairplane maneuverとの言葉の間で、ローマ字に曖昧さは存在し得ない）。また、理解すべきことは、説明を容易にするために、図8に示される辞書は、他の点においても単純化されている。例えば、実施の確率辞書は、それぞれの言葉の、さらに多くのマッピングを含み得、または、所定の確率閾値を超過するマッピングを含み得る。

#### 【0056】

本発明の好ましい実施形態は、そのような確率辞書を用い、一つの言語および/または文字セットにおいて表現されたクエリを、別の言語および/または文字セットに変換し、

10

20

30

40

50

それにより、ユーザが、元々のクエリとは異なる文字セットおよび／または言語にて書かれた文書を発見することを可能にする。例えば、ユーザがローマ字で「cars」というクエリを入力した場合、確率辞書は、「cars」というローマ字の言葉を、「cars」という漢字の言葉へとマッピングするために使用され得る。この方法において、クエリの文字セット（例えばローマ字）とマッチングする文書の文字セット（例えば漢字）と同じではない場合でさえも、ユーザは、それらのクエリに関連する文書を発見することができる。この特定の例において、クエリの実際の言語は変化せず（ローマ字および漢字は日本語を表現するために使用される）、文字エンコードのみが変化することに注意されたい。

#### 【0057】

別の例として、ASCIIの英語における「tired」という言葉は、Latin1の文字エンコードを用いたドイツ語における「müede」という言葉にマッピングされ得る。というのは、ウムラウトuという文字は、ASCIIに存在しないからである。この例において、辞書は、辞書は他の言語に変換され（英語からドイツ語）、他の文字エンコード（ASCIIからLatin1）へと変換されることに注意されたい。

#### 【0058】

好ましい実施形態において、上記されたマッピング辞書は、自動的な方法において構築され、統計的な技術に関連して、ウェブ上で利用可能な情報を用いる。好ましい実施形態は、正確な変換に達するために、異なる言語および／または文字セットにて書かれたアンカーテキストなどのような、パラレルで連係されたバイリンガルのコーパスを用いる。このデータを用いて、好ましい実施形態は、言葉のマッピングの辞書を構築することが可能である。これは、例えば、単に言語 $S_i$ （ソース言語）が、連係されたテキストの対（例えばアンカー、文、など）におけるトークン $T_j$ （ターゲット言語）と同時に生じる回数を数えることによってなされ得る。しかしながら、任意の適切な技術が用いられ得ることは理解されたい。十分に広く正確に連係されたセットが存在しない場合において、この方法は、相対的に曖昧な多対多のマッピングを生成し得る。したがって、例えば、 $S_1$ は、一部の確率を用いて、 $T_2$ 、 $T_3$ 、 $T_7$ および $T_8$ にマッピングされ得ることが決定され得る。しかしながら、これは、以下で詳細に記載するように、許容可能であり、一部の実施形態において、追加的な改良が、それぞれのマッピングの個々の可能性（例えば、以前のユーザのクエリ、結果ページ上のアイテムのユーザ選択および／またはそのようなものを調べることによって）を増加させるためになされ得る。

#### 【0059】

図9は、確率辞書を構築するための、パラレルアンカーテキストの使用を示す。アンカーテキストは、ウェブページ間（または、所定のウェブページ内の位置）のハイパーリンクに関連付けられたテキストを含む。例えば、ハイパーテキストマークアップ言語（HTML）において、「`<A href = "http://www.abc.com">Banks and Savings and Loans</A>`」というコマンドは、「Banks and Savings and Loans」というテキストを、`http://www.abc.com`のウェブページを提示するハイパーリンクとして表示させる。「Banks and Savings and Loans」というテキストは、アンカーテキストと呼ばれ、通常は、提示されるウェブページ（例えば、`www.abc.com`）の短い記載を提供する。実際は、アンカーテキストは、しばしば、そのページ自体よりも、より正確なウェブページの記載を提供し、提示するウェブページの性質を決定するのに、特に有用であり得る。付け加えて、アンカーテキストにおける言葉の使用および配布は、しばしば、趣旨および長さにおいて、ユーザのクエリにて見出されるものに近い。また、所定のページを提示するアンカーの多くは、同じか、または高度に類似したテキストである場合がある。例えば、`www.google.com`を提示するアンカーは、単に、「Google」であるか、または、他のテキストとともに、この言葉を少なくとも使用する。したがって、例えばカタカナなど、`www.google.com`を提示する全てのアンカーを検証することで、「Google」のカタカナ変換は、最も高

10

20

30

40

50

い頻度で現れる言葉を単に探すことによって、相対的に高い信頼度でもって推測され得る（可能であれば、単なる「ここをクリック」といったような、所定の低い情報内容を除去した後で）。本発明の好ましい実施形態は、正確な変換を提供するために、アンカーテキストのこれらの特性の利点を持つ。

#### 【0060】

図9を参照すると、第1の文字セット（例えばASCII）に書かれた言葉を含むクエリを受け取ると（ブロック902）、サーバは、言葉にあるアンカーテキストのセットを検証し得る（ブロック904）。例えば、サーバは、その言葉を含むそれらのアンカーを識別するために、全ての既知のアンカーのインデックスを検証し得る。次に、それらのアンカーが提示するウェブページは識別され（ブロック906）、アンカーは、それらのページを提示するターゲット言語またはターゲット文字セット（例えば、ひらがな、カタカナ、および/または漢字）で書かれる（ブロック908）。システムはここで、2つの文書のセットを有する（そこでは、アンカーテキストは、文書のフォームと見なされる）。一つの文書のセット（例えば、もともとのASCIIのクエリを含むアンカー）におけるクエリという言葉の分配は、次いで、他の文書セット（例えば、パラレルアンカー）における変換されたフレーズの最も適当な候補を識別するために使用される。統計は、アンカーテキストが現れる頻度に関して計算され得、これらの統計は、アンカーテキストに見出される言葉の相対的な頻度または確率が、もともとのクエリの正しい変換であるかどうかを決定するために使用され得る（ブロック910）。多数の言葉を有するクエリに対して、上記のプロセスは、それぞれの言葉に対して繰り返され得、または、全体のクエリは、単一の言葉として扱われ得、あるいは、一部の他の適切な言葉の群が使用され得る。例えば、クエリが「big houses」である場合、可能な変換の辞書は、そのフレーズを含む、連係されたアンカーテキストを見出すことによって構築され得る。同様に、クエリが3つ以上の言葉を含む場合、適切なマッピングを決定するための経験は、クエリという言葉の適切なサブセットを取り上げ、それらの言葉の結果を生成することによって構築され得る。

#### 【0061】

図9に示される方法において変換を実行する利点は、変換システムが一つの言語または文字セットにおける言葉とターゲットセットにおけるそれらとの間のマッピングの予備的知識を必要としない。その代わり、マッピングは、統計的分析を実行するために利用可能であるデータの本体に基づいて、ダイナミックに決定され得る。したがって、例えば、従来のスタティックな辞書を維持する労力または費用（例えば、言語的分析および調査）を負うことなしに、スラング、イディオム、適切な名前などに対する正確な変換を発見することが可能である。

#### 【0062】

前述の変換の例示的な実施形態は、ここで、図10～図12と関連して記載される。この例において、ユーザは「house」というクエリという言葉を入力し、スペイン語で書かれた検索結果（または、単に、クエリという言葉の変換）を得ることを望んでいると想定する。サーバは、英語の「house」を、スペイン語の相当物に変換することを企てる。

#### 【0063】

図10を参照すると、様々なウェブページ959、961、963、965が、アンカーテキスト960、962、964、966を介して、ページ972および974にリンクされる。一部のページおよびそれらに関連するアンカーテキストは、英語で書かれており（つまり、ページ959a～959eおよび963a～963t）、一部はスペイン語で書かれている（ページ961a～961eおよび965a～965j）。サーバは、第1に、「house」という言葉を使用する全てのアンカーの位置を突き止める。これらのアンカーは、例えば、サーバにおいて格納されたアンカーテキストのインデックスを検索することによって、位置が突き止められ得る。そのようなインデックスを使用して、サーバは第1に5つのアンカー960を見出し得、それぞれが「big house」というフレーズを使用し、ウェブページ972を提示する。サーバはまた、次に、ページ97

2を提示する5つのターゲット言語（例えばスペイン語）のアンカー962が存在することを決定する。図10に示される例において、これらのアンカーは「c a s a g r a n d e」というテキストを含む。同じページ（アンカー960およびアンカー962など）、またはそれに対して所定の関係を有するページに提示されるアンカーは、「連係（a l i g n e d）」されているといい、そこでは、より一般的な意味では、配置が通常、連係されたアイテムの相当物（またはほぼ相当物）を意味する。

#### 【0064】

図11Aは、それぞれのターゲット言語の言葉は、ターゲット言語のアンカー962に現れる頻度を示す。図11Aに示されるように、「c a s a」および「g r a n d e」はそれぞれ、5回現れる（つまり、それぞれのアンカー962に一度）。したがって、ターゲットアンカー962に現れる、トータルで10の言葉（つまり、5つのアンカーのそれぞれにおいて、アンカー毎に2つの言葉）から、「c a s a」は半分を占め、「g r a n d e」はもう半分を占める。したがって、図11Aで示されるように、この時点で、「h o u s e」という言葉は、「c a s a」および「g r a n d e」の両方の言葉が等しい頻度で現れるゆえに、等しい確率で、「c a s a」または「g r a n d e」のいずれかにマッピングされ得る。

#### 【0065】

しかしながら、図10にて示されるように、システムはまた、「h o u s e」という言葉を含む20個の英語のアンカー964を見出し、ページ974に提示し、ならびに、「c a s a」という言葉を含む10個のスペイン語のアンカー966を見出し、ページ974を提示する。図11Bにて示されるように、「h o u s e」という言葉は、「c a s a」という言葉に、0.75の確率（つまり、 $15/20$ ）でマッピングされ、「g r a n d e」という言葉に、0.25の確率（つまり $5/20$ ）でマッピングされる。これらの確率は、そのターゲット言語のアンカーにおけるそれぞれの言葉の出現回数のトータル（「c a s a」の場合は15）を、ターゲット言語のアンカーにおける言葉のトータル数（重複を含む）によって、単に割ることによって計算される（つまり、20の言葉とは、アンカー962に含まれる10、および964に含まれる10）。代替的には、または追加的には、他の技術が、所定の変換またはマッピングの確率を計算および/または改良するために使用され得る。例えば、当業者が理解するように、任意の様々な既知の技術は、ベイズ法（B a y e s i a n m e t h o d s）、ヒストグラムスムージング（h i s t o g r a m s m o o t h i n g）、カーネルスムージング（k e r n e l s m o o t h i n g）、縮小推定量（s h r i n k a g e e s t i m a t o r s）、および/または他の推定方法などの確率推定の分散エラーを減少させるために使用され得る。

#### 【0066】

さらなるアンカーテキストが利用可能である場合、確率は、一層さらに改良され得る。例えば、最終確率分配は、図12にて示されたものと同様であり得、「h o u s e」は、相対的に高確率で、「c a s a」、および、それに接尾語が付いた形式「c a s i t a」にマッピングされ、幾分か低い確率で、「c a s i n o」および

#### 【0067】

#### 【数1】

“mansión”

（スペイン語でm a n s i o nに相当する）にマッピングされ、僅かな確率で、「g r a n d e」にマッピングされる。したがって、正確な変換、およびほぼ同意語の識別は、変換される言語および/または文字セットの知識なしで得られ得る。

#### 【0068】

クエリの変換すると、サーバは、変換を使用して検索を実行し得る。例えば、ユーザは「h o t e l s i n K y o t o」に対するローマ字のクエリを入力する場合、上記された技術は、サーバに、カタカナ、ひらがな、および漢字のクエリのフォームを推測させることができ、それらのクエリを使用して検索を実行させることができ、ならびに

10

20

30

40

50

、適切なユーザインターフェース内で、ユーザへのそれらのクエリのフォームのそれぞれに対する、組み合わせられた結果を提示させることができる。

【 0 0 6 9 】

図 1 0 ~ 図 1 2 に関連して記載された例は、例示のために提供されるのであり、限定のためではなく、多くの変更がそこに表された方法論に対してなされ得ることは理解されるべきである。例えば、異なる統計的な技術が一定の確率に達するために使用され得、および/または、修正は、上記された基本的な技術に対してなされ得る。同様に、上記された変換技術は、単に、ユーザによって入力された言葉またはフレーズの変換を実行するために使用され得、関連するインターネット検索を実行し、または確率辞書を作成するために使用される必要はない。付け加えて、先行する例は、ユーザのクエリの受け取りのアプリケーションとに生じるものとして、変換プロセスを記載するが、他の実施形態においては、マッピングプロセスは、ユーザのクエリが受け取られる前に実行され得ることは理解されるべきである。そのような計算前のマッピングは、図 8 に記載されたような辞書に格納され得、次いで、それらが受け取られたときに、ユーザのクエリを変換するように適合される。最後に、関係されたアンカーテキストとは異なるテキストは、変換を実行するために使用され得ることは理解されるべきである。例えば、関係された文章または他のデータは、同様な方法にて使用され得る。多くの国では、一つ以上の公式言語または認められた言語が存在しており、新聞や定期刊行物はしばしばこれらの言語のそれぞれにて書かれた同じ記事を含む。これらのパラレルな変換は、言葉の変換の確率辞書を準備するために、前記されたアンカーテキストと同様の方法において使用され得る。

【 0 0 7 0 】

したがって、好ましい実施形態により、有利にも、ユーザは、従来の方法において、検索クエリおよび/または変換要求を入力することができ、正確で自動の変換および検索を提供する。一部の実施形態において、追加の改良が上記された基本的なモデルに対してなされ得る。例えば、一部の実施形態において、優先（加重）は、もともとのクエリおよび/または他の関係されたアンカーにおける言葉の数と同様の言葉の数を含むアンカーに与えられ得る。例えば、図 1 0 に示されるシステムにおいて、優先は、ページ 9 7 4 に提示されるアンカーに与えられ得る。というのは、もともとのクエリと同様に、それらは、それぞれ単一の言葉を含むからである。同様に、「l a c a s a g r a n d e」のテキストを含むアンカーがまたページ 9 7 2 に提示された場合、その加重は、適切な要因によって軽減され得、というのは、それが関係された他のアンカーよりも多くの言葉（例えば 3 つ）を含むからである。そのような加重スキームは、適切な要因によってこれらのアンカーの言葉と関連される頻度を増加させることによって、図 1 1 B に示される確率計算に反映され得る。

【 0 0 7 1 】

上記された変換プロセスはまた、検索それ自体の効率を改善するために利用され得る。例えば、確率辞書は、様々な変換およびもともとのクエリという言葉の同意語などを含む、進行中のクエリを拡張するために使用され得る。文書検索に先立つユーザのクエリを拡張することによって、同じ「概念」の同時検索が実行され得、それにより、検索結果は、ユーザが探しているものを含む可能性を増加させる。代替的には、または追加的には、確率辞書は、文書という言葉の拡張を提供することによって、通常の文書インデックス付けのプロセスを補うために使用され得る。例えば、文書にて見出される言葉は、確率辞書からの変換を用いて、文書のインデックスにおいて補われ得、その文書は、もともとの文書にて見出された同じ言葉を正確に使用しない検索によってさえも、位置を突き止められる可能性を増加させる。

【 0 0 7 2 】

上記された変換技術を使用する場合に生じ得る問題は、データの希薄（例えば、「c a s a」を「h o u s e」にマップすることを最終的に決定するには十分なアンカーがない）、または、多様性の欠如（全てのアンカーが同じものを示す）などのためであり、システムは、十分に正確な確率のマッピングに達することが不可能であり得る。したがって、



一部の実施形態において、確率マッピングは、ユーザの行動を検証することでさらに改良され得る。いくつかの例示的な技術が以下に記載される。

【 0 0 7 3 】

例えば、再度、サーバが「house」に対する変換を得ることを望んでいる場合を想定する。しかしながら、見出され得るアンカーテキストが、「big house」というフレーズ、または「casa grande」というフレーズのいずれかを含むことを想定する。そのアンカーテキストにおける多様性の欠如のために、確率辞書は以下のマッピングに達し得る。

house casa、0.5の確率

house grande、0.5の確率

10

big casa、0.5の確率

big grande、0.5の確率

grande house、0.5の確率

grande big、0.5の確率

casa house、0.5の確率

casa big、0.5の確率

20

ここでユーザが「casa」という言葉を用いて検索エンジンにクエリすると想定する。この時点で、検索エンジンは、「casa」という言葉を含むページを返答し得、また、「house」という言葉をちょうど含むN結果と、「big」という言葉をちょうど含むM結果とを合わせる。実際には、NおよびMは、マッピングに内在する確率を考慮するように調整され得、その結果、相対的に見込みのないマッピングは、より少ない結果が表示されることになる。ユーザが、「big」という言葉を含む結果をクリックするよりも、「house」という言葉を含む結果を10倍、クリックした場合、マッピングの確率は、例えば、以下のように調整される。

house casa、0.9の確率

house grande、0.1の確率

30

big casa、0.1の確率

big grande、0.9の確率

grande house、0.1の確率

grande big、0.9の確率

casa house、0.9の確率

casa big、0.1の確率

40

実際の数値は、クリックが考慮されるユーザの数、その言葉の両方を含むページのクリックの回数、結果のセットの中の、当該の言葉を含む結果の置換、および/またはそれらのようなものなど、他の様々な要因に依存し得ることに注意されたい。この例（つまり、0.1および0.9）に与えられた調整された確率は、例示の目的のためであることは理解されたい。当業者は、上記されたものなどのユーザフィードバックに与えられた実際の加重は、任意の適切な方法にてインプリメントされ得ることは理解する。

【 0 0 7 4 】

また、前述の例は、ユーザフィードバックの使用の例を容易にするように簡易化されていることに注意されたい。例えば、一部のシステムにおいて、所定の変換を実行することを補助するために、他の変換から得られた情報を利用することが可能である。例えば、今

50

提示された例において、「house」という言葉が「big house」という言葉を含んだアンカーテキストに現れる場合、「house」は、「grande」にマッピングされるよりも、より適切に「casa」にマッピングされることを決定することは可能であり得る。例えば、既に、「big」が非常に高確率で、および、十分に大きなデータのセットにわたって、「grande」にマッピングされるように決定されている場合（および、アンカーテキストが同意語のリストで構成されてない場合）、次いで、house-to-casaのマッピングは、たとえ、「house」または「casa」を含むアンカーが要領を得ない場合であっても、house-to-grandeのマッピングを介して、優先を与えられ得る。

【0075】

変換の正確さおよび/または検索結果の有用性はまた、ユーザのクエリセッションの履歴を検証することによって改善され得る。例えば、多くの場合、システムは、ユーザが入力した以前のクエリを知っている（例えば、クッキーまたはサーバ上のユーザアカウントに格納された情報などを介して）。この履歴データは、そのユーザからのクエリの、可能な意味をランクするために使用され得、飛行に関連するものから、フィッシングに関するクエリに対して「bank」を明確にする。したがって、このプロセスは、可能な変換のセットを狭めるために使用され得る。一部の実施形態において、ユーザインターフェースにて、「Xの検索を意図しますか？」（ここで、Xは所定の変換の優先を意味する）などのメッセージに関連して、それらを表示することによって、これらを示唆し得、その一方で、結果の第1のページにおいて、可能な再公式化のそれぞれからの結果の一部を表示する。ユーザが「...意図しますか？」によって示唆された代替のうちの一つを選択するか、結果ページに提示された結果のうちの一つを選択する場合、システムは、クエリという言葉の適当な変換、およびユーザの適当な検索バイアスに関する追加的な証拠を得る。これらの信号の両方は、次いで、システムによって利用され得、一般的な場合、およびユーザに特殊な場合の両方において、言葉のマッピングの適当な根拠を更新する（例えば、確率辞書において）。

【0076】

#### D. 結論

上記されたように、本発明と合致する方法およびシステムは、曖昧な検索結果に応じて検索結果を提供し、言葉を他の文字セットおよび/または言語に変換するために使用され得る。様々な変換および検索技術、ならびにシステムが記載されている。しかしながら、前述の記載は、例示のために提示されるものであり、多くの修正および変更が、上記の教示に照らし合わせて、または、本発明の実施を通して、可能であることは理解されたい。例えば、前述の記載はクライアント-サーバ構成に基づいているが、当業者は、ピアツーピア構成もまた、本発明に合致して使用され得ることは理解されたい。さらに、記載された実施はソフトウェアを含むが、本発明は、ハードウェアとソフトウェアとの組み合わせ、またはハードウェアのみとして実施され得る。付け加えて、本発明の局面は、メモリに格納されるように記載されたが、当業者は、これらの局面もまた、ハードディスク、フロッピーディスク（登録商標）、またはCD-ROMなどのような二次的格納装置、インターネットからの搬送波、あるいは、RAMまたはROMの他の形式などの、他のタイプのコンピュータ可読媒体に格納され得る。本発明の範囲は、それゆえ、請求項およびそれらの均等物によって定義される。

【図面の簡単な説明】

【0077】

【図1】本発明と合致する方法および装置がインプリメントされ得るシステムのブロック図を例示する。

【図2】本発明と合致する、クライアントデバイスのブロック図を例示する。

【図3】3つの文書を例示する図である。

【図4a】従来の英数字のインデックスを例示する。

【図4b】従来の英数字の検索クエリに応じて、検索結果を提供するフロー図である。

【図 5 a】曖昧な検索クエリに応じた検索結果を提供するための、本発明に合致したフロー図を例示する。

【図 5 b】数字の情報を数字の情報にマッピングする図を例示する。

【図 5 c】（記載なし）

【図 6】曖昧な検索クエリに応じた検索結果を提供するための、本発明に合致した別のフロー図を例示する。

【図 7】本発明の実施形態に従い、検索を実行する方法を例示する。

【図 8】文字セットの変換の確率辞書を例示する。

【図 9】確率辞書を構築するためのパラレルアンカーテキストの使用を例示する。

【図 10】アンカーテキストを使用してリンクされた文書の集まりを例示する。

【図 11】図 11 A および図 11 B は、図 10 に示されたアンカーテキストに基づく適当な変換の計算を例示する。

【図 12】例示的な言葉の変換と関連した確率分配を示す。

10

【図 1】

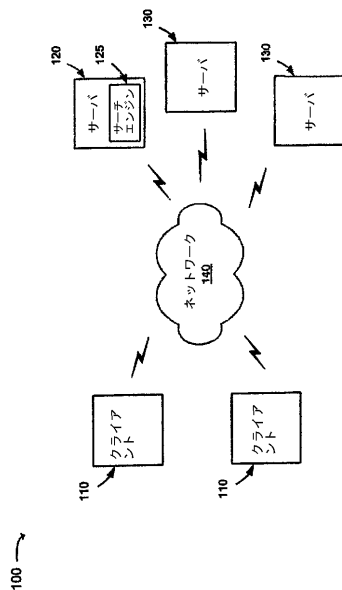


FIG. 1

【図 2】

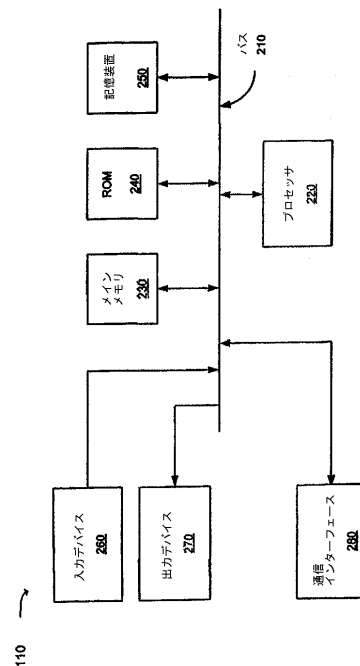


FIG. 2

【図 3】

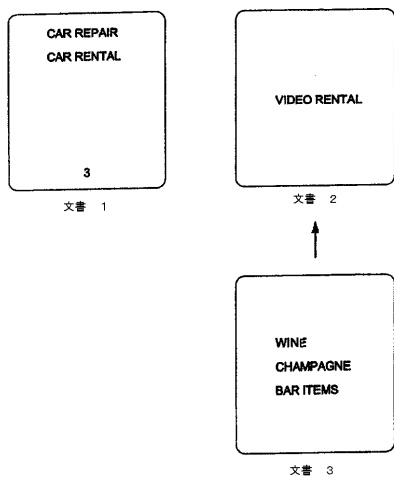


FIG. 3

【図 4 a】

言葉	位置 (文書)
3	文書 1
BAR	文書 3
CAR	文書 1
CHAMPAGNE	文書 3
ITEMS	文書 3
RENTAL	文書 1 および 2
REPAIR	文書 1
VIDEO	文書 2
WINE	文書 3

FIG. 4A

【図 4 b】

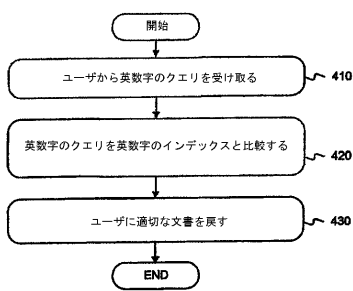


FIG. 4B

【図 5 a】

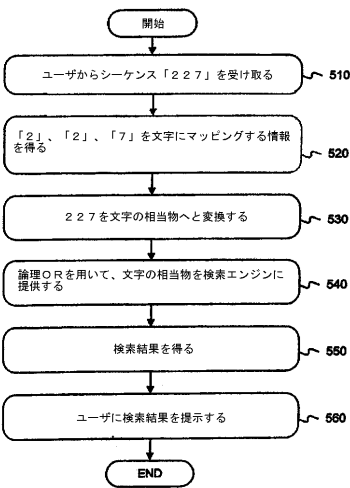
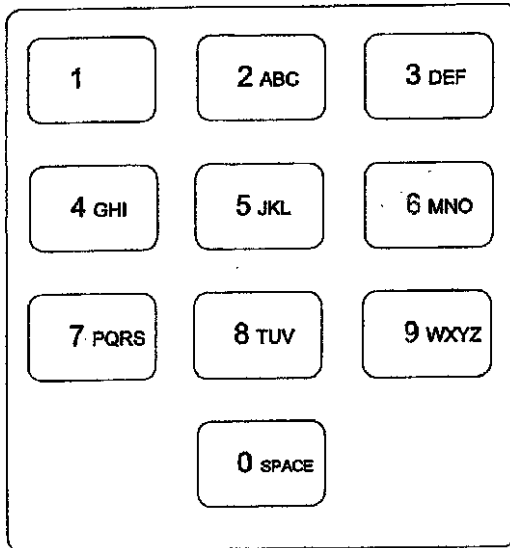


FIG. 5A

【図 5 B】



【図 5 c】

言葉	位置 (文書)
3	文書 1
227	文書 1 および 3
242872463	文書 3
48367	文書 3
736825	文書 1 および 2
737247	文書 1
84336	文書 2
8463	文書 3

FIG. 5C

FIG. 5B

【図 6】

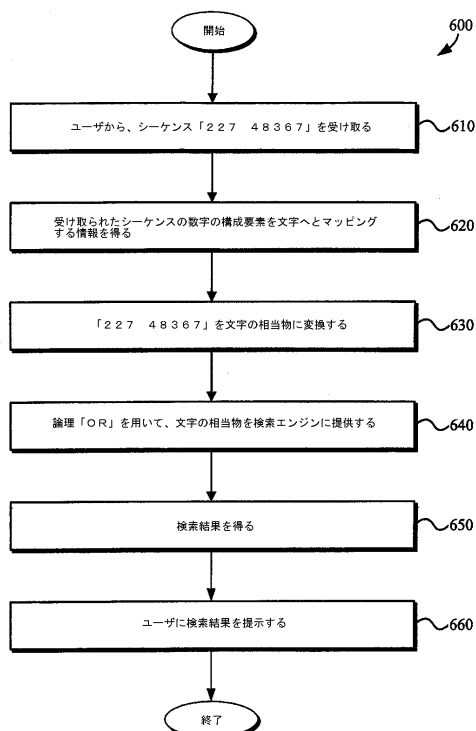


FIG. 6

【図 7】

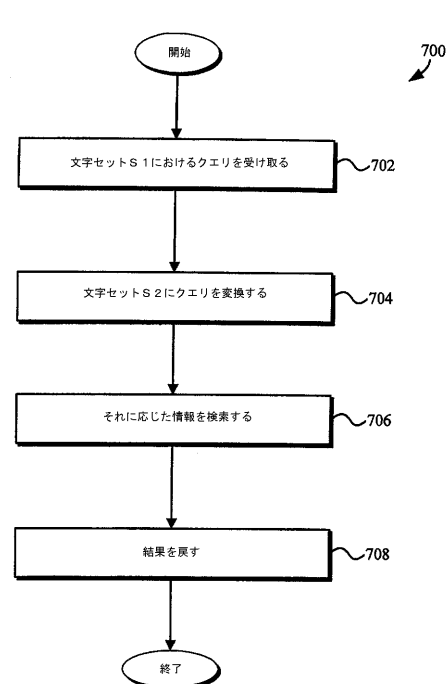


FIG. 7

【図 8】

808	802 ローマ字の言葉	804 漢字の言葉	806 確率 (%)
808	<Bank> romaji	<Financial institution> kanji	0.4
		<Steep slope> kanji	0.3
		<Airplane maneuver> kanji	0.2
		<Other> kanji	0.1
810	<Car> romaji	<Automobile> kanji	0.9
		<Other> kanji	0.1
812	<House> romaji	<A dwelling> kanji	0.7
		<To contain> kanji	0.25
814	<Plane> romaji	<Airplane> kanji	0.6
		<Flat surface> kanji	0.25
		<Carpenter's tool> kanji	0.1
		<Other> kanji	0.05

FIG. 8

【図 9】

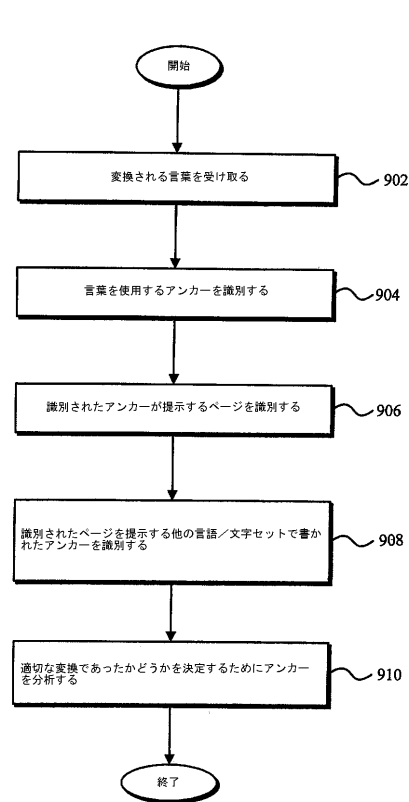


FIG. 9

【図 10】

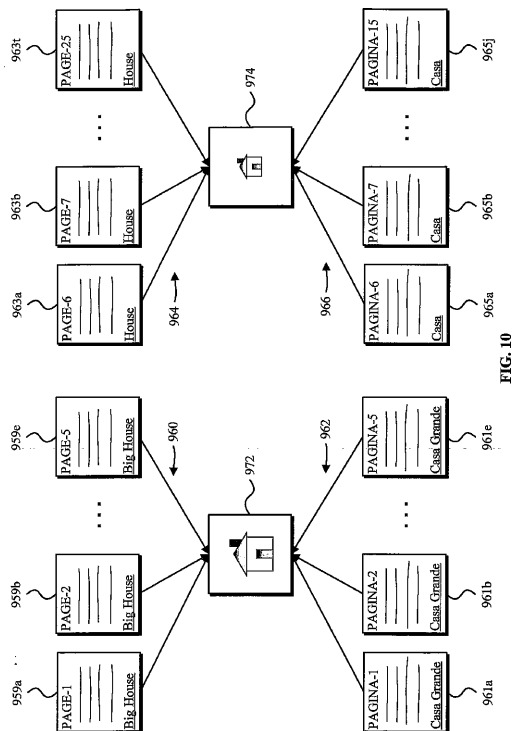


FIG. 10

【図 11】

House	1	Casa	$\# = (\text{house} \times \text{casa}): 5/10 = 0.5$
	2	Casa	
	3	Casa	
	4	Casa	
	5	Casa	$\# = (\text{house} \times \text{grande}): 5/10 = 0.5$
	6	Grande	
	7	Grande	
	8	Grande	
	9	Grande	
	10	Grande	

FIG. 11A

House	1	Casa	$\# = (\text{house} \times \text{casa}): 15/20 = 0.75$
	2	Casa	
	3	Casa	
	...	Casa	
	13	Casa	$\# = (\text{house} \times \text{grande}): 5/20 = 0.25$
	14	Casa	
	15	Casa	
	16	Grande	
	17	Grande	
	18	Grande	
	19	Grande	
	20	Grande	

FIG. 11B

【図 12】

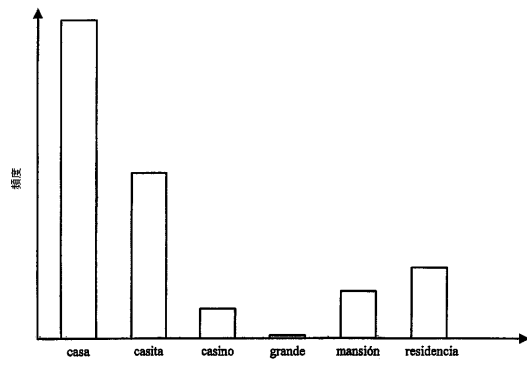


FIG. 12

## フロントページの続き

- (72)発明者 ボンテ, ジェイ エム.  
アメリカ合衆国 カリフォルニア 94043, マウンテン ビュー, マージョリー コート  
2439
- (72)発明者 サハミ, メヘラーン  
アメリカ合衆国 カリフォルニア 94063, レッドウッド シティ, フーバー ストリー  
ト 3238
- (72)発明者 ゲマワット, サンジャイ  
アメリカ合衆国 カリフォルニア 94043, マウンテン ビュー, ノース レンフストル  
フ アベニュー 111, ナンバー184
- (72)発明者 バウアー, ジョン エー.  
アメリカ合衆国 カリフォルニア 94040, マウンテン ビュー, デル メディオ アベ  
ニュー 415, ナンバー8

審査官 波内 みさ

- (56)参考文献 特開2000-163441(JP, A)  
特開2002-024266(JP, A)  
国際公開第2003/058374(WO, A1)  
特開2002-222189(JP, A)  
特開2002-259374(JP, A)  
菊井 玄一郎, 言語の壁を越えて文書を検索する - クロスランゲージ情報検索 -, 人工知能学会  
誌 第15巻 第4号, 日本, 社団法人人工知能学会, 2000年 7月 1日, 第15巻 第4  
号, 550~558  
乾 健太郎, 言語表現を言い換える技術, 言語処理学会第8回年次大会チュートリアル資料, 日  
本, 言語処理学会, 2002年 3月17日, 1~21

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

G06F 17/28