

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4668567号  
(P4668567)

(45) 発行日 平成23年4月13日(2011.4.13)

(24) 登録日 平成23年1月21日(2011.1.21)

(51) Int. Cl.		F I			
<b>G06F 13/00</b>	<b>(2006.01)</b>		G06F 13/00	540A	
<b>G06F 15/00</b>	<b>(2006.01)</b>		G06F 15/00	310U	
<b>G06F 17/30</b>	<b>(2006.01)</b>		G06F 17/30	180Z	

請求項の数 20 外国語出願 (全 28 頁)

(21) 出願番号	特願2004-239997 (P2004-239997)	(73) 特許権者	500046438
(22) 出願日	平成16年8月19日 (2004.8.19)		マイクロソフト コーポレーション
(65) 公開番号	特開2005-135381 (P2005-135381A)		アメリカ合衆国 ワシントン州 9805
(43) 公開日	平成17年5月26日 (2005.5.26)		2-6399 レッドモンド ワン マイ
審査請求日	平成19年8月20日 (2007.8.20)		クロソフト ウェイ
(31) 優先権主張番号	10/670, 681	(74) 代理人	100077481
(32) 優先日	平成15年9月25日 (2003.9.25)		弁理士 谷 義一
(33) 優先権主張国	米国 (US)	(74) 代理人	100088915
			弁理士 阿部 和夫
		(72) 発明者	エリック ディー. プリル
			アメリカ合衆国 98052 ワシントン
			州 レッドモンド 172 コート ノー
			スイースト 3042

最終頁に続く

(54) 【発明の名称】 クライアントベースのウェブクロウリングのためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項1】

ウェブクローラと前記ウェブクローラによって収集されたウェブページ情報に関する第1のデータセットを格納した第1のストレージとを含む第1のコンピュータと、

ブラウザまたはプロキシサーバから訪れたウェブページのウェブページ情報に関する第2のデータセットを格納した第2のストレージを含む少なくとも1つの第2のコンピュータと、

を備え、前記第1のコンピュータは前記第1のデータセットから第1のウェブページに関する第1のウェブページ情報の第1の表現を生成して、その第1の表現を前記第2のコンピュータに送信し、前記第2のコンピュータは前記第1の表現に対応する前記第1のウェブページのウェブページ情報をブラウザまたはプロキシサーバを用いて収集し、前記収集したウェブページ情報の第2の表現が前記第1の表現と異なる場合、前記第2のコンピュータは収集したウェブページ情報に基づき前記第2のデータセットを更新して、その更新された第2のデータセットを前記第1のコンピュータに送信し、前記第1のコンピュータは前記第2のデータセットに基づき第1のデータセットを更新することを特徴とするデータ分析システム。

【請求項2】

ウェブクローラと前記ウェブクローラによって収集されたウェブページ情報に関する第1のデータセットを格納した第1のストレージとを含む第1のコンピュータと、

ブラウザまたはプロキシサーバから訪れたウェブページのウェブページ情報に関する第

2 のデータセットを格納した第 2 のストレージを含む少なくとも 1 つの第 2 のコンピュータと、

を備え、前記第 1 のコンピュータは前記第 1 のデータセットから第 1 のウェブページに関する第 1 のウェブページ情報の第 1 の表現を生成して、その第 1 の表現を前記第 2 のコンピュータに送信し、前記第 2 のコンピュータは前記第 1 の表現に対応する前記第 1 のウェブページのウェブページ情報をブラウザまたはプロキシサーバを用いて収集し、前記第 2 のコンピュータは収集したウェブページ情報に基づき前記第 2 のデータセットを更新して、その更新された第 2 のデータセットを前記第 1 のコンピュータに送信し、前記第 1 のコンピュータは、前記第 2 のコンピュータから受信したウェブページ情報の第 2 の表現が前記第 1 の表現と異なる場合、前記第 2 のデータセットに基づき第 1 のデータセットを更新することを特徴とするデータ分析システム。

10

**【請求項 3】**

前記ウェブクローラは、インターネットウェブクローラを含むことを特徴とする請求項 1 又は 2 に記載のシステム。

**【請求項 4】**

前記ウェブクローラは、イントラネットウェブクローラを含むことを特徴とする請求項 1 又は 2 に記載のシステム。

**【請求項 5】**

前記第 1 のコンピュータは、前記少なくとも 1 つの第 2 のコンピュータからの前記第 2 のデータセットの受信を制御するためのスケジューリング機能を提供することを特徴とする請求項 1 又は 2 に記載のシステム。

20

**【請求項 6】**

前記第 2 のコンピュータは、前記第 1 のデータセットと前記第 2 のデータセットを比較して、前記第 1 のコンピュータのウェブクローラによって検索されたなりすましデータを検出するのにさらに使用されることを特徴とする請求項 1 又は 2 に記載のシステム。

**【請求項 7】**

前記第 2 のコンピュータは、前記第 1 のデータセットに関連するデータについての状況情報を生成するのにさらに使用され、前記状況情報は、前記第 2 のコンピュータが複数ある場合、少なくとも 1 つの他の第 2 のコンピュータに送信されることを特徴とする請求項 1 又は 2 に記載のシステム。

30

**【請求項 8】**

前記状況情報は、前記第 1 のデータセットに関連する情報の鮮度を示す鮮度フラグを少なくとも一部は含むことを特徴とする請求項 7 に記載のシステム。

**【請求項 9】**

前記状況情報は、前記第 1 のデータセットに関連する情報のコンテンツのハッシュを少なくとも一部は含むことを特徴とする請求項 7 に記載のシステム。

**【請求項 10】**

前記状況情報は、前記第 1 のデータセットに関連する情報のコピーを少なくとも一部は含むことを特徴とする請求項 7 に記載のシステム。

**【請求項 11】**

第 1 のコンピュータが備えるウェブクローラによってウェブページ情報に関する第 1 のデータセットを生成して前記第 1 のコンピュータが備える第 1 のストレージに前記第 1 のデータセットを格納するステップと、

40

前記第 1 のコンピュータが第 1 のウェブページに関する第 1 のウェブページ情報の第 1 の表現を前記第 1 のデータセットに基づき生成するステップと、

前記第 1 のコンピュータからブラウザまたはプロキシサーバから訪れたウェブページのウェブページ情報に関する第 2 のデータセットを格納した第 2 のストレージを含む少なくとも 1 つの第 2 のコンピュータに前記第 1 の表現を送信するステップと、

第 2 のコンピュータが前記第 1 の表現に対応する前記第 1 のウェブページのウェブページ情報を収集するステップと、

50

第2のコンピュータが前記第1の表現に基づき収集したウェブページ情報の第2の表現を生成するステップと、

第2のコンピュータが、前記第2の表現が前記第1の表現と異なる場合、前記収集したウェブページ情報に基づき第2のデータセットを更新して、その更新された前記第2のデータセットを前記第1のコンピュータに送信するステップと、

前記第1のコンピュータが受信した前記前記第2のデータセットに基づき前記第1のデータセットを更新するステップと

を有することを特徴とするデータ分析方法。

【請求項12】

第1のコンピュータが備えるウェブクローラによってウェブページ情報に関する第1のデータセットを生成して前記第1のコンピュータが備える第1のストレージに前記第1のデータセットを格納するステップと、

前記第1のコンピュータが第1のウェブページに関する第1のウェブページ情報の第1の表現を前記第1のデータセットに基づき生成するステップと、

前記第1のコンピュータからブラウザまたはプロキシサーバから訪れたウェブページのウェブページ情報に関する第2のデータセットを格納した第2のストレージを含む少なくとも1つの第2のコンピュータに前記第1の表現を送信するステップと、

第2のコンピュータが前記第1の表現に対応する前記第1のウェブページのウェブページ情報を収集するステップと、

第2のコンピュータが、前記収集したウェブページ情報に基づき第2のデータセットを更新して、その更新された前記第2のデータセットを前記第1のコンピュータに送信するステップと、

前記第1のコンピュータが前記第2のデータセットに基づき、前記第1のウェブページに関するウェブページ情報の第2の表現を生成するステップと、

前記第1のコンピュータが、前記第2の表現が前記第1の表現と異なる場合、受信した前記前記第2のデータセットに基づき前記第1のデータセットを更新するステップと

を有することを特徴とするデータ分析方法。

【請求項13】

前記ウェブクローラは、インターネットウェブクローラを含むことを特徴とする請求項11又は12に記載の方法。

【請求項14】

前記ウェブクローラは、イントラネットウェブクローラを含むことを特徴とする請求項11又は12に記載の方法。

【請求項15】

前記第1のコンピュータが、前記少なくとも1つの第2のコンピュータからの前記第2のデータセットの受信を制御するスケジューリングするステップをさらに有することを特徴とする方法。

【請求項16】

前記第2のコンピュータが、前記第1のデータセットと前記第2のデータセットを比較して、前記第1のコンピュータのウェブクローラによって検索されたなりすましデータを検出するステップをさらに有することを特徴とする請求項11又は12に記載の方法。

【請求項17】

前記第2のコンピュータが、前記第1のデータセットに関連するデータについての状況情報を生成するステップをさらに有し、前記状況情報は、前記第2のコンピュータが複数ある場合、少なくとも1つの他の第2のコンピュータに送信されることを特徴とする請求項11又は12に記載の方法。

【請求項18】

前記状況情報は、前記第1のデータセットに関連する情報の鮮度を示す鮮度フラグを少なくとも一部は含むことを特徴とする請求項17に記載の方法。

【請求項19】

10

20

30

40

50

前記状況情報は、前記第1のデータセットに関連する情報のコンテンツのハッシュを少なくとも一部は含むことを特徴とする請求項17に記載の方法。

【請求項20】

前記状況情報は、前記第1のデータセットに関連する情報のコピーを少なくとも一部は含むことを特徴とする請求項17に記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、概してデータ分析に関し、より詳細には、分散型ウェブクローラを使って、ネットワーク接続されたシステムから情報を取得するシステムおよび方法に関する。

10

【背景技術】

【0002】

高コスト、低性能のデータ処理システムから、低コスト、高性能の通信システム、問題解決システム、および娯楽システムへの、コンピュータおよびネットワーク技術の発展により、書簡のやり取り、請求書の支払い、買物、予算の立案、および情報収集など、日常業務を実施するための負担を軽減する、コスト効率が高く時間を節約する手段がもたらされた。たとえば、有線または無線技術を介してインターネットとインターフェイスをとる計算機システムは、世界中に位置するウェブサイトおよびサーバのリポジトリからの大量の情報に、ユーザが指一本で、ほぼ瞬時に近いアクセスをするためのチャンネルをユーザに提供する。

20

【0003】

一般に、ウェブサイトおよびサーバを介して利用可能な情報は、ウェブクライアント（たとえばコンピュータ）上で実行されるウェブブラウザを介してアクセスされる。たとえば、ウェブユーザは、ウェブブラウザを展開し、ウェブサイトのURL（Uniform Resource Locator）（たとえば、ウェブアドレスおよび/またはインターネットアドレス）をウェブブラウザのアドレスバーに入力し、キーボード上のエンターキーを押下しまたはマウスで「go」ボタンをクリックすることによって、ウェブサイトにアクセスすることができる。URLは通常、アクセスを容易にする4つの情報を含む。すなわち、情報交換のための規則および標準の集合を示すプロトコル（互いに通信するためのコンピュータ用語）と、ウェブサイトまでの位置指定と、ウェブサイトを維持する組織名と、組織のタイプを識別する添字（たとえば、com、org、net、gov、およびedu）である。

30

【0004】

場合によっては、ユーザは、自分がアクセスしたいと望むサイトもしくはサーバの名称、および/またはサイトもしくはサーバへのURLを事前に知っている。このような状況において、ユーザは、上述したように、アドレスバーにURLを入力しサイトに接続することによって、サイトにアクセスすることができる。しかし、ほとんどの場合、ユーザは、URLもサイト名も知らない。代わりに、ユーザは、検索エンジンを利用して、自分が提供したキーワードに基づいてサイトの発見を容易にする。概して、検索エンジンは、キーワードを求めてウェブサイトおよびサーバのコンテンツを検索するとともに、キーワードが見つかったウェブサイトおよびサーバへのリンクの一覧を返す、実行可能なアプリケーションまたはプログラムからなる。基本的に、検索エンジンは、関連づけられたURLとしてできるだけ多くの文書を検索するウェブ「クローラ（crawler）」（別名、「スパイダー」または「ロボット」）を組み込む。この情報は次いで、インデクサ（indexer）が、検索されたデータを処理することができるように格納される。インデクサは、文書を読み出し、各文書に含まれるキーワードおよび文書の他の属性に基づいて、優先順位をつけられた索引を作成する。それぞれの検索エンジンは一般に、所有権のあるアルゴリズムを利用して、クエリに対して有意義な結果が返されるように、索引を作成する。

40

【0005】

50

したがって、ウェブクローラは、検索エンジンの動作にとって重要である。現在および最新の検索結果を提供するために、クローラは、ウェブを絶えず検索して、新しいウェブページを見つけ、古いウェブページ情報をアップデートし、消去されたページを削除しなければならない。インターネット上で見られるウェブページの数はいくつもの天文学的である。したがって、ウェブクローラは、極度に高速であることが要求される。ほとんどのウェブクローラは、ウェブページを提供するサーバにポーリングを行うことによってデータを集めるので、クローラは、ある特定のサーバにアクセスするとき、できるだけ控えめにもしなければならない。そうでないと、クローラは、サーバの資源すべてを非常に高速に吸収し、サーバをシャットダウンさせてしまう場合がある。一般に、クローラは、サーバのウェブページにアクセスする前に、サーバに対してそれ自体を識別し、アクセス許可を求める。この時点で、サーバは、サーバの資源すべてを盗む不正クローラに対して、アクセスを拒否することができる。サーバをホスティングするウェブページは一般に、検索エンジンにより、ユーザがウェブページをより容易にみつけることが可能になるので、検索エンジンから利益を受ける。したがって、ほとんどのサーバは、サーバの資源すべてを使い果たさない限り、クローラを歓迎し、そうすることによって、サーバのコンテンツは、ユーザによってより便利に活用することができるようになる。

10

【発明の開示】

【発明が解決しようとする課題】

【0006】

サーバに対してそれ自体を識別するクローラの欠点の1つは、サーバがクローラに対して「なりすまし(spoof)」を行うことができることである。サーバは通常、広範なインターネットに対して公開したくない、保護された領域を有する。クローラがそれ自体を識別するとき、クローラは、どの領域にアクセスすることができないかも知られる。クローラは、その特定のサーバとの仕事上の関係を維持したい場合、サーバの要求を遵守する。しかし、サーバは、その本当のコンテンツについてなりすましを行いまは偽りたい場合、そのサーバの本当のURLを模倣しているが「代替」コンテンツを含むページ領域にクローラを向けさせることができる。したがって、猫に関する情報のみを通常は提供するサーバが、ウェブクローラのみがアクセスすることができるセクション中に、犬に関する情報を有するURLを設定することができる。これは、ユーザが「犬」を検索するとき、猫に関するサーバのウェブページが検索エンジンによって示されるように行われる。通常、なりすましは、サーバのコンテンツが世間によっていかにわしいとみなされているが、サーバが、通常の「キーワード」の範囲を超えてそのコンテンツを広めたいと思う場合に使用される。このようにして、いかにわしい素材が、たとえば花、犬、猫、天気など、一般的な言葉を使うことによって、検索エンジンの一覧で返される場合がある。なりすましは、正確さも、なりすましを受けたウェブクローラのデータを利用する検索エンジンの評判も低下させる。

20

30

【課題を解決するための手段】

【0007】

以下では、本発明のいくつかの態様の基本的な理解をもたらすために、本発明の簡略な要約を提示する。この要約は、本発明の包括的な概要ではない。本発明の主要な/重大な要素を明らかにすることも、本発明の範囲を詳述することも意図していない。後で提示するより詳細な説明の前置きとして、本発明のいくつかの概念を簡略な形で提示することだけを目的としている。

40

【0008】

本発明は、一般にデータ分析に関し、より詳細には、分散型ウェブクローラを使って、ネットワーク接続されたシステムから情報を取得するシステムおよび方法に関する。サーバのクライアントの、分散される性質は、高速かつ正確なウェブクローリングデータを提供するために利用される。サーバのウェブクローラによって集められた情報は、クローラのデータをアップデートするために、サーバのクライアントによって検索されたデータと比較される。本発明の一例では、データの比較は、検索エンジンの結果ページを介して広

50

められる情報を使うことによって遂行される。本発明の別の例では、データの妥当性確認は、ウェブクロウラのデータを要約する、サーバから生じる、クライアントの辞書によって遂行される。本発明の一態様では、「弱インジケータ(weak indicator)」関数の集合の、ある弱インジケータ関数が、クライアントにランダムに送られる。こうした弱インジケータ関数は、サーバのウェブクロウラによって見つけれられるすべてのURLの全一覧より著しく小さく、したがって、サーバとクライアントの間の通信トラフィックを大幅に削減する。このことは、サーバとクライアントの間のインターフェイスの簡素化を容易にするとともに、ウェブクロウラのデータの正確さを最適に保つ。

【発明の効果】

【0009】

本発明は、ウェブクロウラがなりすましを受けないよう阻止する手段を提供し、データの正確さを向上することによって、データ分析も容易にする。本発明を利用するサーバは、ウェブクロウラのデータを、クライアントによって提供されるデータと比較することによって、なりすましを阻止することができる。このことは、サーバが、その検索エンジンから、なりすましたデータを除去することを可能にし、より高品質の検索エンジン結果を可能にする。この能力は、特に害のない検索において通常は返されない、いかがわしい素材をフィルタリングして取り除くことを容易にし、検索エンジンのクライアントに対して、よりユーザフレンドリーな体験をもたらす。

【0010】

上記の目的および関連する目的を達成するために、本発明の例示的な態様を、本明細書において、以下の記述および添付の図面に関連して説明する。ただし、こうした態様は本発明の原理を利用することができる様々な方法のごくわずかを示すに過ぎず、本発明は、このようなすべての態様およびその等価物を含むことを意図する。本発明の他の利点および新規の特徴は、本発明の以下の詳細な説明を図面と併せ読むことにより、明らかになるであろう。

【発明を実施するための最良の形態】

【0011】

ここで図面を参照して本発明を説明するが、同じ参照番号は、全体を通して同じ要素を指すのに使われる。以下の記述では、説明のために、多くの具体的な詳細を、本発明の完全な理解をもたらすために述べる。ただし、こうした具体的な詳細なしでも本発明を実施

【0012】

本願において使う「コンポーネント」という用語は、コンピュータ関連のエンティティ、すなわちハードウェア、ハードウェアおよびソフトウェアの組合せ、ソフトウェア、または実行中のソフトウェアのいずれかを指すことを意図している。たとえば、コンポーネントは、プロセッサ上で実行されている処理、プロセッサ、オブジェクト、実行ファイル、実行スレッド、プログラム、および/またはコンピュータでよいが、それに限定されない。例として、サーバ上で実行されているアプリケーションおよびそのサーバ両方がコンピュータコンポーネントとなることができる。1つまたは複数のコンポーネントが実行の処理をし、および/またはスレッド中に常駐することができ、コンポーネントは、1台のコンピュータに配置することも、および/または2台以上のコンピュータの間に分散することもできる。「スレッド」とは、オペレーティングシステムのカーネルが実行のためにスケジューリングする処理におけるエンティティである。当該分野において公知であるように、各スレッドは、スレッドの実行に関連づけられた揮発性データである、関連する「コンテキスト」を有する。スレッドのコンテキストは、システムレジスタのコンテンツおよびスレッドの処理に属する仮想アドレスを含む。したがって、スレッドのコンテキストを含む実際のデータは、実行時に変化する。

【0013】

本発明は、ウェブ文書の索引を維持する、改良されたシステムおよび方法を提供する。

索引は、他のタイプの情報のデータを検索し、維持するのにも利用されることができる。従来のウェブクローラは、本発明によって軽減されるある欠点を有する。各クライアント（たとえば、ウェブにアクセスする任意の人のマシン）は、ローカルな情報を格納するので、したがって、クライアントが最後に訪れたときから、ウェブページが変更されているか否かを知ることができる。変更されている場合、クライアントは、この情報を検索エンジンに伝達することができる。同様に、サーバは、クライアントが訪れたウェブページについての情報を用いて、現時点でサーバにとって未知であるページを見つけることができる。文書を効率よく見つけ、そうした文書についての現時点での知識を維持することは、イントラネットおよびインターネット検索両方にとって、非常に重要なタスクである。本発明は、イントラネット検索などの状況でも利用することができ、その場合、ページをク

10

## 【0014】

検索エンジンの、（インターネット、イントラネット、またはそれ以外にとって）重要なコンポーネントは、データクローラまたは文書クローラである。文書クローラは、2つの主要なタスクを実施する。すなわち、検索エンジンによって索引を付けられるべき未知の文書を見つけないこと、および、その文書が、既知の各文書についての最新の知識を有することを保証しようと試みることである。こうしたタスクは両方とも、困難であり、（ページランクの質とともに）検索エンジンにおいて最も重要であり、目に見える品質の差別化要因に属する。文書クローラは一般に、サーバモデルに基づく。検索エンジンは、トポロジ検索によってウェブをクローラする。既知のウェブページからなるシードセット（se

20

## 【0015】

現在の、サーバベースのクローリングパラダイムには、いくつかの脆弱性がある。第1に、クローラは、シード文書の1つで始まる後続リンクによって到達することができるページしか発見することができない。最近の研究によると、大部分のウェブページは、どの検索エンジンによっても現時点で索引を付けられていないことがわかっている。第2に、検索エンジンは、クローラがページを再訪したときに、文書に対する変更（たとえば、コンテンツの変更や、既に存在しないページ）を知ることしかできない。

30

## 【0016】

本発明は、上述した脆弱性を改善するやり方で、文書（たとえばデータ）を効率的に見つけ、既知の文書についての最新の知識を保持するシステムおよび方法を提供すること、分散型の、クライアントベースのクローラによって達成される。各クライアント（たとえば、ウェブにアクセスする任意の人のマシン）は、ローカルな情報を格納し、したがって、クライアントが最後に訪れたときから、ウェブページが変更されているかどうかを知ることができる。変更されている場合、クライアントは、この情報を検索エンジンに伝達することができる。同様に、サーバは、クライアントが訪れたウェブページについての情報を用いて、現時点でサーバにとって未知であるページを見つけることができる。

40

## 【0017】

図1に、本発明の一態様によるデータ分析システム100のブロック図を示してある。本発明のこの事例において、データ分析システム100は、1から「N」（Nは、1から無限大のどの数も表す）の番号がつけられたクライアント102～106、通信システム108、検索サーバ110、およびウェブページサーバ112からなる。クライアント102～106は、検索サーバ110用のウェブページ情報の「分散資源」群を備える。こうしたクライアントは一般に、新しいURL、およびウェブページの変更などを、通信シ

50

システム108を介して検索サーバ110に提供するように機能する。通信システム108は、インターネットおよび/またはイントラネットなどからなる。通信システム108は、検索サーバ110とクライアント102~106の間の通信用アクセス手段を提供する。通信システム108は、ウェブページ情報を集めるために、クライアント102~106と他のウェブページサーバ112、および/または検索サーバ110と他のサーバの間の通信も可能にする。本質的に、ウェブクロウラの機能性は、検索サーバ内で機能するだけでなく、検索サーバ110およびクライアント102~106に分散される。検索サーバ110は、クライアント102~106を使って、ウェブページサーバ112から情報を取得し、検索サーバ自体の情報の洗練を容易にする。この機能を分散させることによって、本発明は、検索エンジンがそのデータを使用することができる、より最新の、頑強な、なりすましを防ぐデータセットを提供する。

10

**【0018】**

図2に移ると、本発明の態様によるデータ分析システム200の別のブロック図を示してある。データ分析システム200は、クライアント202およびサーバ204からなり、それらの間で相互運用可能な通信手段を有する。通常の動作の間、サーバ204は、ウェブクロウラをホスティングし、そのウェブクロウラは、ウェブページをホスティングする他のサーバを求めて、インターネットなどの通信ネットワークを検索する。クロウラは、ウェブページの検索エンジンでの使用のために、こうしたウェブページについての情報ソースを編集する。サーバ204は次いで、このウェブページ情報の表現をクライアント202に送る。この表現により、クライアント202は、その特定のウェブページをホスティングするサーバにアクセスするとき、独立してウェブページ情報を検証することが可能になる。クライアント202は、サーバ204にとって未知であるウェブページを検出することもできる。これにより、クライアント202は、既知および未知のウェブページについての、変更/状況および/または新しい情報を編集することが可能になる。この情報は次いで、サーバ204に送信される。サーバ204は、この情報を使って、サーバの、クロウラ側の元のウェブページデータを改良する。分散資源を有することにより、サーバ204は、それ自体の直接資源(たとえば、プロセッサの使用、記憶空間など)に負荷をかけることなく、そのクロウラ機能を拡張する。さらに、ウェブクロウラは通常、アクセスする各サーバに対してそれ自体を識別するので、そのサーバ上の誤ったデータに宛先変更されるという危険がある。サーバも、ウェブクロウラがサーバの資源に対して課し得るアクセス量および時間を制限することができる。サーバにアクセスするクライアントは一般に、こうした制限をもたず、誤ったデータに宛先変更されない。したがって、クライアントのウェブページデータは、ウェブクロウラによって編集された、誤ったデータを修正するのに使うことができる。本発明のこの態様は、後でより詳しく説明する。

20

30

**【0019】**

図3を参照すると、本発明の態様によるデータ分析システム300のさらに別のブロック図を示してある。データ分析システム300は、クライアントシステムコンポーネント302およびサーバシステムコンポーネント304からなり、その間で相互運用可能な通信システム(CS)306を有する。本発明のこの事例において、クライアントシステムコンポーネント302は、CSインターフェイスコンポーネント308、クライアント制御コンポーネント310、およびデータ記憶コンポーネント312からなる。CS GUI(グラフィカルユーザインターフェイス)コンポーネント308は、通常は、利用される通信システムのタイプに特有のインターフェイスをユーザに提供する。このようなインターフェイスの一例は、少なくともワールドワイドウェブ上の情報を図表によって中継するために使われるウェブブラウザである。ウェブブラウザは、一企業内で供給されるウェブページなどのイントラネットを「サーフィンする」のにも使うことができる。本発明の他の例では、同様の情報を、グラフィカルユーザインターフェイスではなく、テキストベースのインターフェイスなどを使って中継することができる。一般に、ユーザが、通信システム306に接続されたりリモートサーバ上に常駐する検索エンジンにおいて検索クエリを実行することを可能にするのは、このコンポーネント308である。したがって、CS

40

50

GUIコンポーネント308は、通信システム306から情報を受信し、かつ/または送信する。クライアント制御コンポーネント310は、ウェブクロールを容易にすることに關するクライアントの制御を提供する。クライアント制御コンポーネント310は、たとえばウェブページなどの情報に關するデータを受信し、かつ/または送信する。クライアント制御コンポーネント310は、アルゴリズムを処理し、データの変更および状況を追跡し、かつ/またはデータ分析システム300内のクライアント用のローカルデータ記憶を制御する。コンポーネント310は、CS GUIコンポーネント308からの情報を、ウェブクローラから受け取った情報を用いて分析して、違いなどを判別することもできる。クライアント制御コンポーネント310は、クライアントが、ウェブクローラなどのための「分散資源」として関与することも可能にする。コンポーネント310は、格納されたデータにアクセスし、CS GUIコンポーネント308に情報を提供することもできる。本発明の一例では、CS GUIコンポーネント308は、埋め込まれたクローラのデータを送信し、かつ/または受信する。したがって、クライアント制御コンポーネント310は、GUIコンポーネント308とインターフェイスをとって、埋め込まれたクローラ関連データを、このようにして受信し、かつ/または送信する。同様に、コンポーネント310は、サーバからの制御を、同じやり方で送信し、かつ/または受信することもできる。本発明の別の例では、クライアント制御コンポーネント310は、サーバと同様に振る舞い、ピアツーピア方式で他のクライアントに制御を提供することができる。クライアント制御コンポーネント310およびCS GUIコンポーネントの機能は、単一のコンポーネント内で組み合わせることができることが、当業者には理解されよう。クライアントを、CS GUIコンポーネント308をもたない分散資源として使うことも可能である。本発明のこの場合の例は、別のクライアントを中継し、かつ/または制御するクライアントを含むが、それに限定されない。データ記憶コンポーネント312は、たとえば、サーバからのクローラのデータ、クライアントからのクローラのデータ、ウェブページの変更、新しいウェブページデータ、およびクライアント制御パラメータなどを格納するのに使われる。コンポーネント312は、利用される本発明の例に応じて、クライアント制御コンポーネント310および/またはCS GUIコンポーネント308と直接、インターフェイスをとることができる。データ記憶コンポーネント312は、ハードドライブ、ランダムアクセスメモリ、読出し専用メモリ、取外し可能媒体、およびCD-ROMなどのデータ記憶装置でもよい。本発明のさらに別の例では、データ記憶コンポーネント312に格納された情報は、CS GUIコンポーネント308ともクライアント制御コンポーネント310ともインターフェイスをとることなく、サーバによって直接アクセスすることができる。いくつかの例では、このことは、より高速なデータの検索を可能にする。

#### 【0020】

本発明の一例では、通信システム306は、「インターネット」などの相互接続ネットワークでよい。通信システム306は、WAN（ワイドエリアネットワーク）および/またはLAN（ローカルエリアネットワーク）などのイントラネットシステムでもよい。通信システム306は、より従来型の通信手段、たとえば、電話システム、無線システム、光信号（光学）システム、およびサウンドシステムなどを利用することもできる。他のグローバルおよびローカルネットワーク構造を、本発明によって通信システム306として使うこともできることが当業者には理解されよう。

#### 【0021】

サーバシステムコンポーネント304は、検索エンジンコンポーネント314、分散資源制御コンポーネント316、クローラコンポーネント318、データ記憶コンポーネント320、および任意選択によるCSデータホスティングコンポーネント322からなる。本発明の一例では、クローラコンポーネント318は、サーバおよび/またはプロキシサーバにアクセスするのに通信システム306を使って、ウェブページ関連情報、たとえばウェブページのコンテンツ、古さ、サイズ、URL、および埋込みリンクなどを取得する。この情報は次いで、データ記憶コンポーネント320に格納される。データ記憶コン

10

20

30

40

50

ポーネント 3 2 0 は、ハードドライブ、ランダムアクセスメモリ、読出し専用メモリ、取  
外し可能媒体、および CD-ROM などのデータ記憶装置でよい。検索エンジンコンポー  
ネント 3 1 4 は、ウェブクローラ 3 1 8 によって示されるとともにデータ記憶コンポー  
ネント 3 2 0 に格納されるすべてのウェブページ用の検索機能を提供する。検索エンジン  
コンポーネント 3 1 4 は、ユーザから検索要求 / クエリを受け取り、データ記憶コンポー  
ネント 3 2 0 上の情報にアクセスしてリンク一覧およびウェブページデータを編集して、ユ  
ーザに送信する。したがって、一般的なシステムでは、検索コンポーネント 3 1 4 は、ク  
ローラコンポーネント 3 1 8 によって取得される情報にのみ依拠することができる。しか  
し、本発明の例では、分散資源コントローラ 3 1 6 は、データ記憶コンポーネント 3 2 0  
に格納された情報の編集を容易にし、情報がより頑強、最新、より包括的になるようにす  
る。分散資源制御コンポーネント 3 1 6 は、たとえば、単一の分散型クローラ、すなわち  
「クライアントベースのウェブクローラ」として相互作用する、サーバのクライアントな  
どの分散資源に対する制御を提供する。コンポーネント 3 1 6 は、クライアントシステム  
コンポーネント 3 0 2 などの分散資源から受信したデータの分析、機能ならびにデータの  
割振りと割振りのタイミングの判定、既知のクローラのデータを判定するための分散資源  
へのアルゴリズムの提供、データのアップデートおよび / または追加の受信、データのア  
ップデートおよび / または追加のデータ記憶コンポーネント 3 2 0 への格納、分散資源の  
最適化された利用の決定、ある特定の検索クエリに対する検索結果ページへのデータの埋  
込みを可能にするための検索エンジンコンポーネント 3 1 4 へのページデータの提供、埋  
込みページリンク情報を含むページを生成するためのページデータのインターネットのサ  
ービスプロバイダへの提供、カウント、タイプ、なりすましの割合、およびソースなどの  
データ特性の追跡などの機能を提供する。本発明の別の例では、コンポーネント 3 1 6 が  
通信システム 3 0 6 に直接アクセスするのではなく、検索ページコンポーネント 3 1 4 が  
、分散資源制御コンポーネント 3 1 6 用の情報を送信し、かつ / または受信する。

10

20

#### 【 0 0 2 2 】

本発明の例では、任意選択の CS データホスティングコンポーネント 3 2 2 は、通信シ  
ステム 3 0 6 および分散資源制御コンポーネント 3 1 6 両方とインターフェイスをとる。  
CS データホスティングコンポーネント 3 2 2 は、ウェブページのホスト機能を提供し、  
ユーザにウェブページへのアクセスを提供する。CS データホスティングコンポーネント  
3 2 2 は、分散資源制御コンポーネント 3 1 6 と対話するので、ウェブページのリンク情  
報を受信し、コンポーネント 3 2 2 がホスティングするウェブページに情報を直接埋め込  
むことができる。本発明の他の例では、CS データホスティングコンポーネント 3 2 2 は  
、データ記憶コンポーネント 3 2 0 と直接インターフェイスをとって、ウェブページに埋  
め込むための情報にアクセスする。本発明のさらに別の例では、CS データホスティング  
コンポーネント 3 2 2 は、検索エンジンコンポーネント 3 1 4 とインターフェイスをとっ  
て、そのウェブページのリンクに埋め込むための情報にアクセスする。本発明のさらに別  
の例では、CS データホスティングコンポーネント 3 2 2 は、クライアントなどの分散資  
源に常駐することができる。コンポーネント 3 2 2 は、サーバシステムコンポーネント 3  
0 4 へのアクセス権を有する別のサーバに常駐することもできる。この例では、クライ  
アント (またはサーバ) は、事実上、ホスティングされるウェブページに対するサーバとな  
り、ウェブページのリンクに埋め込むための情報を、そのローカルストレージおよび / ま  
たは他のローカル手段から供給する。

30

40

#### 【 0 0 2 3 】

各コンポーネントをそれぞれ独立に説明したが、本発明の他の例におけるコンポーネ  
ントは、他のコンポーネントに関連づけられた機能を含むことができることが当業者には理  
解されよう。同様に、いくつかのコンポーネントは、本発明の範囲を変えることなく、削  
除することができる。

#### 【 0 0 2 4 】

図 4 に移ると、本発明の態様によるデータ分析システム 4 0 0 を示すさらに別のプロッ  
ク図を示してある。データ分析システム 4 0 0 は、クライアントシステムコンポーネント

50

402およびサーバシステムコンポーネント404からなり、その間で相互運用可能な通信システム406を有する。本発明のこの事例において、サーバシステムコンポーネント404は、分散資源制御コンポーネント414およびデータ記憶コンポーネント416からなる。サーバシステムコンポーネント404は、クライアントシステムコンポーネント402からのウェブページ情報の受信に関して、本発明の事例を強調するために省略してある。通常、情報は、通信システム406を介して、分散資源制御コンポーネント414へ、およびコンポーネント414から流れる。クライアントシステムコンポーネント402は、クライアント制御コンポーネント408、データ記憶コンポーネント410、および任意選択の通知コンポーネント412からなる。本発明のこの事例において、通知コンポーネント412は、クライアントシステムコンポーネント402からサーバシステムコンポーネント404に流れるデータを制御する。本発明の他の例では、コンポーネント412は、クライアントシステムコンポーネント402と他のクライアントシステムコンポーネントとの間のピアツーピア通信も制御する。具体的には、通知コンポーネント412は、いつ、および/またはどのデータが、クライアントシステムコンポーネント402から送信されるべきかを決定する。決定は、蓄積されたウェブページのデータのサイズ、サーバシステムコンポーネント404にとって未知であるリンクが見つかったかどうか、ウェブページに対する変更の重要度(たとえば50%以上のコンテンツの変更および/または重要度の高いページの変更など)、アクセス許可時刻、および/または分散資源制御コンポーネント414によって設定される一般的なアクセス許可時間などに基づくことができる。通知コンポーネント412は、データ転送のために、アルゴリズムを用いて、独自の重要度の要素および/または独自のタイミングスケジュールを決定することもできる。通知コンポーネント412の機能は、クライアント制御コンポーネント408および/または図4に示さない他のクライアントシステムコンポーネントに常駐できることが、当業者には理解されよう。

#### 【0025】

本発明の完全な理解のために、動作例を説明する。本発明の一例では、分散型クライアントベースのクロウラは、以下のように動作する。潜在的な新しいウェブページ、およびウェブページのコンテンツ/状況変更についての着信クライアントメッセージを受信するサーバ、ならびにサーバと通信するクライアントの集合が存在するものと仮定する。クライアントマシンは、ウェブブラウジング用に使われるパーソナルコンピュータ、またはパーソナルコンピュータにページを供給するプロキシサーバのいずれでもよい。クライアントは、(1)ウェブページに到達するのに使われるURL、(2)ウェブページのコンテンツのハッシュ、(3)ウェブページのコンテンツ、および(4)訪れた時間を含むことができるがそれに限定されない、閲覧されるウェブページ上の情報を集めるように装備される。本発明のいくつかの例(たとえばプロキシサーバなど)において、この情報すべてを保持することは実現不可能であり、いくつかの情報は、ある程度の期間保持されるだけである。

#### 【0026】

本発明の別の例では、クライアントは、ある特定のブラウザまたはプロキシサーバから訪れたウェブページのURLを、一定の期間記録し、次いで、このURLの集合をサーバに送信する。サーバは次いで、どのURLがサーバにとって未知であったかを調べ、そうしたURLを、今後のクロール/ダウンロード/索引づけのために、既知のURL一覧に追加する。こうすることにより、サーバに関連づけられた検索エンジンは、トポロジカルなクロールによって見つけることができなかつたウェブページについて知ることができるようになる。

#### 【0027】

クライアントからサーバに送信される情報のボリュームを減らすために、クライアントは、ある特定のURLをすでにサーバに知らせてあるか否かという情報をローカルに保持することができる。まだ知らせていない場合は、サーバに情報を送信するだけでよい。2つのウェブページが同じであるかどうかを効率的に判定する公知の方法がある。ハッシュ関

10

20

30

40

50

数によって、各文書を整数にマッピングし、次いで、2つのハッシュ値が同じであることを調べる。URLに関連づけられたコンテンツの最新のハッシュが、そのURLに関連づけられたコンテンツの、以前のハッシュと異なる場合、そのコンテンツは変更されている。クライアントは、ウェブページを訪れる度に、そのページのハッシュ値を計算する。クライアントは、そのページを訪れたことがある場合、ハッシュ値が変わっているかを調べる。ハッシュ値が変わっている場合、クライアントは、クライアントが最後にそのページを訪れた後にウェブページが変更されたと判定し、サーバに知らせることができる。クライアントは、新しい<url、ハッシュ値>のペアをローカルに記録する。

**【0028】**

クライアントが、変更についてサーバに知らせるための、異なるいくつかの方法がある。最も簡単な方法は、URLのコンテンツ/状況が変更されたというメッセージを単に送信するだけである。次いで、サーバは、そのページをできるだけすぐに再クロールするようにスケジュールすることができる。サーバがページを再訪する必要をなくするために、クライアントは、付加情報を送信することができる。クライアントは、最後に訪ずれたときの、ページをキャッシュしたコピーをもっている場合、古いバージョンと新しいバージョンの間の違いとともに、古いハッシュ値、および新しいハッシュ値を送信することができる。サーバは最初に、クライアントの古いハッシュ値が、そのページの、サーバの現在のハッシュ値と一致するかを調べる。一致する場合、サーバは、それに従ってページのコンテンツをアップデートすることができる。一部の文書変更は、他の変更より重要であることに留意されたい。たとえば、ある場合には、ページ全体が変更されているが、別の場合には、ただ1つのコンマがある文に追加されているだけである。クライアントは、変更の重要度を計算し、(a)この情報を使って、どのアップデートをサーバに送信するかという優先順位を決定するか、または(b)他のページ情報とともに重要度の値をサーバに送信し、そうすることによって、サーバがページの再クロール/再索引づけに優先順位をつける際にこの情報を利用できるようにすることができる。変更重要度関数の例は、変更された文書の割合、変更の言語的/意味的重要度、および変更によって影響を受けるユーザ検索の割合の推定などのような項目を含むが、それに限定されない。重要度は、ページの人気の推定によって重みづけすることもできる。

**【0029】**

上述した通信手段に伴う欠点の1つは、クライアントとサーバの間の、重大なトラフィックのオーバーヘッドを生じることである。たとえば、100個のクライアントがすべて、ページ「X」を初めて訪れる場合、クライアントはそれぞれ、ページ「X」を発見したというメッセージをサーバに送信する。同様に、サーバは、ページ「Y」が変更されたことを通知されると、それ以外のクライアントからその通知を受ける必要はない。したがって、クライアントとサーバの間の不必要な通信を大幅に減少させる、本発明のこれ以外の例を後で説明する。

**【0030】**

図5を参照すると、本発明の態様による、ページ検索結果を使用するデータ分析システム500の図を示してある。データ分析システム500は、検索結果ページ506を有するクライアント502と、クライアント502からサーバ504へ送信し(508)受信する(510)ための通信手段を有する検索サーバ504とからなる。本発明の事例の第1の例において、クライアント502は、サーバ504に、変更されたウェブページを通知するが、いかなる付加情報も送信しない。ユーザが検索エンジンを使用すると、検索サーバ504は、検索結果ページ506中の各ウェブページに関して、コンテンツの、サーバ側バージョンのハッシュと、コンテンツが、新しくないと知られているか否かを示す鮮度フラグとを含む結果とを、クライアントに提供する。クライアント502は、検索結果ページ506にあるページの1つを訪れる場合、最初に、サーバ504が、ページが新しくないことを知っているか否かを調べ(たとえば、別のクライアントがサーバ504に知らせたが、サーバ504がそのページを更新していない)、ページのコンテンツのハッシュを計算し、検索エンジンが提供したハッシュと比較する。2つのハッシュが一致しない

10

20

30

40

50

場合、クライアント502は、そのURLに関連づけられたコンテンツが変更されたという通知をサーバ504に送る。サーバ504は、通知を受け取ると、鮮度フラグの状況を変更し、再クロールのために、優先待ち行列にそのURLを追加する。

#### 【0031】

この例は、クライアントが、サーバにページ差異情報（この情報は、サーバが、ウェブページについてのサーバ側の情報を、クロールせずにアップデートするため、および/またはサーバがウェブページをいつ再クロールすべきかという優先順位をつけるために使うことができる）を送るシナリオ用に拡張することができる。この拡張は、検索エンジンに、各検索結果を有する2つの追加フィールド、すなわち最新クライアント通知の時間、および最新クライアント通知からのページのハッシュ値を送信させることによって、遂行することができる。クライアントが、検索エンジンによって返されたページを訪れて、(a) known-not-freshフラグが偽であるか、または(b) known-not-freshフラグが真であり、かつ最新クライアント通知からのハッシュ値が、このクライアントがページに対して計算したハッシュ値と異なる場合、クライアントは、サーバに通知を行う。ページ変更の周期を認識し、そうすることによって、ページがAからB、C、Aへと繰り返し変更される場合、本発明によりその変更を認識し、このページについてのクライアントによるアップデートを制限できるようにすることも可能である。

10

#### 【0032】

クライアントとサーバの間の不必要な通信の量の低下に加え、「検索エンジンの結果ページによるメッセージ通信」のそれ以外の利点の1つは、サーバが既に知っているウェブページについての情報のみをクライアントがサーバに送ることを保証することによって、秘密に関する起こり得るいくつかの問題を回避することである。このようにして、クライアントは、たとえば、クライアントが秘密にしておくことを期待しているページを訪れず、サーバにこのページの存在を知らせないことが保証される。

20

#### 【0033】

上述した本発明の例に伴う欠点の1つは、サーバが、ユーザの検索クエリを介してクライアントに返したウェブページについての情報しか知ることができないことである。この要件は、サーバに、どの検索エンジンを介してクライアントに返されたウェブページについても知らせることによって、緩和することができる。クライアントは、ユーザがいかなる検索エンジンを使っていることも認識するように装備される。クライアントは、検索結果を訪ねるとき、コンテンツのハッシュを計算する。クライアントは、このURLを訪れたことがある場合、コンテンツのハッシュをキャッシュしている。2つのハッシュが異なる場合、クライアントは、URLおよび新しいハッシュを（前回の訪問からの経過時間、および他の情報に応じて）サーバにアップロードすることができる。クライアントは、そのURLを訪れたことがない場合、URLおよび新しいハッシュをサーバにアップロードすることができる。

30

#### 【0034】

しかし、クライアントは、サーバにとって既知であるURLの詳細な一覧のローカルなコピーをもっている場合、新規なものである可能性のあるURLに遭遇すると、そのURLが既知のURL一覧にあるかを単に調べ、一覧にない場合は、そのURLをサーバに送るだけである。同様に、クライアントは、サーバにとって既知であるすべてのURLに対して、<url、ハッシュ値>のペアの完全な一覧のローカルなコピーをもっているとすると、情報がサーバにとって新しいものである場合にアップデート情報を送るだけでよい。このアイデアに伴う問題は、こうした一覧全体を各クライアントに渡すのは実現不可能であることである。たとえば、検索エンジンは、数ギガバイトものデータとなる何十億ものURLについて知っている場合がある。重大な帯域幅の問題に加え、各クライアントがこのような一覧のためにこれ程大量のローカルストレージを費やすことを期待するのは現実的でない。

40

#### 【0035】

あるいは、本発明の別の例では、重大な帯域幅の問題を排除する通信手段が提供される

50

。たとえば、アルファベット  $S$  が与えられると仮定する。この場合、 $S^*$  は、 $S$  以降の文字からなる全文字列の集合である。辞書  $D$  を、集合  $S^*$  中の文字列の部分集合であると定義する。辞書  $D$  用のインジケータ関数  $I$ 、すなわち  $I : S^* \rightarrow \{0, 1\}$  は、 $d \in D$  ( $d \notin D$ ) である場合、かつその場合に限り、 $I(d) = 1$  というプロパティを有する。辞書  $D$  用の弱インジケータ関数  $I_w$  は、 $d$  が  $D$  中にあることを意味する  $I_w(d) = 1$  (言い換えると、すべての  $d \in D$  に対して  $I_w(d) = 1$  であり、 $I_w(d)$  は、 $D$  中にある任意の  $d$  に対して、0 または 1 のいずれでもよい) というプロパティを有する関数である。最後に、弱インジケータ関数の適切な集合  $I = \{I_{w1}, I_{w2}, \dots, I_{wn}\}$  を、 $D$  中にある任意の  $d$  に対して、 $I_{wi}(d) = 0$  であるような、少なくとも 1 つの  $I_{wi} \in I$  が存在するというプロパティを有する弱インジケータ関数の有限集合であると定義する。

10

## 【0036】

したがって、各クライアントは、 $I$  からランダムに選ばれた弱インジケータ関数を受け取る。こうしたインジケータ関数は、URL の集合全体より大幅に小さく、したがって、そうしたインジケータ関数をクライアントに送ることは現実的である。サーバによって知られているどの URL に対しても、インジケータ関数は、URL が既知であると正しく判定する。サーバによって知られていない URL に対して、インジケータ関数は、既知であると誤ってラベル付けする可能性があるが、その場合、クライアントは、何も知らないか、または未知であると正しくラベルづけし、この場合、クライアントはサーバに知らせることができる。弱インジケータ関数の適切な集合の定義により、サーバにとって未知であるウェブサイトをクライアントによって訪問されるときはいつでも、クライアントのインジケータ関数とそのサイトを新しいものであると認識する確率がゼロでないことが保証される。

20

## 【0037】

上で挙げた例をさらに簡略化するために、 $S = \{a, b, c, d\}$  であり、 $S^*$  中のすべての文字列は長さが 4 未満であり、辞書  $D = \{abcd, adcb, bcdb, ddd\}$  であると仮定する。この辞書用の弱インジケータ関数の例は、以下のようになる。

## 【0038】

(第 2 の文字が  $\{b, d, null\}$  の 1 つ) である場合、かつその場合に限り、 $I(\text{文字列}) = 1$

30

弱インジケータ関数は、 $D$  に対して以下のようにランダムに構成することができる。

## 【0039】

(1)  $D$  を、2 つの重ならない部分辞書  $D'$  および  $D''$  にランダムに区切る。

(2) 「 $i$  番目の文字が集合  $S$  のメンバーである ( $S$  は、 $S$  の部分集合である)」という様式の 1 つまたは複数の項の結合からなる弱インジケータ関数  $I'$  を、 $D'$  に対してランダムに選ぶ。

(3) 同じようにして、 $D''$  用の弱インジケータ関数  $I''$  をランダムに選ぶ。

(4)  $I'(x) = 1$  または  $I''(x) = 1$  である場合、かつその場合に限り、関数  $I(x) = 1$  を作成する。

## 【0040】

40

このようなすべての弱インジケータ関数の集合は、弱インジケータ関数の適切な集合を生じる。クライアント辞書は、 $\langle \text{url}, \text{ページのハッシュ値} \rangle$  のペアからなる辞書をもつことによって、ページの最新性の検出という問題にも拡張することができる。

## 【0041】

本発明の独自の一態様は、専用クローラの視点およびクライアントの視点から、クローラのデータを比較できることである。このことは、サーバの高度化が進む際、特に重要である。「より精密な」ソフトウェアを用いることによって、サーバは、サーバ中にあるデータの流れおよびアクセスをよりうまく制御することができる。これは、任意のまたはすべてのユーザが、サーバ上にある情報の一部または全てにアクセスするのを阻止できることを含む。異なるタイプのユーザに対して、サーバアクセス特権、さらにアクセス時間特

50

権に関して異なる「許可レベル」を与えることもできる。概して、こうした柔軟性の向上は、セキュリティ、有料アクセスの実施、および悪意のあるハッキングの防止など、建設的な目的のために利用される。しかし、サーバ上にあるウェブページの実コンテンツをマスクするのに利用されることも多い。図6に、本発明の態様による、ウェブクロラシステム602を伴うなりすまし処理600のブロック図を示してある。処理600は、ウェブクロラシステム602およびサーバ604を含み、その間で相互運用可能な通信システム606を有する。ウェブクロラシステム602は、クロラコンポーネント608およびデータ記憶コンポーネント610からなる。サーバ604は、サーバアクセス制御612、なりすましデータ614、および実データ616からなる。通常のクロラコンポーネント608がサーバ604にアクセスすると、コンポーネント608は、サーバ604に対してそれ自体をウェブクロラとして識別する。この識別は、「礼儀正しい」とみなされる。礼儀正しさは、サーバ規則を無視することによってサーバを悪用するウェブクロラが、一般に、将来にわたってサーバへのアクセスを拒否されるという点において、自己検閲である。サーバへのアクセスを拒否されることは、サーバアクセスに頼り、検索エンジンのユーザにコンテンツを提供する検索エンジンにとって、特に危機的である。したがって、クロラは通常、礼儀規則を遵守する。他の礼儀規則は、時間限定アクセス、サーバ資源の使用、およびデータの非破壊的な検索などを含む。この例では、サーバアクセス制御612は、クロラコンポーネント608を識別し、実データ616にアクセスを向けるのではなく、クロラコンポーネント608をなりすましデータ614に向ける。なりすましデータは通常、実データ616と同じURL情報を含むが、異なるコンテンツを有する。これは、一般に、いかがわしいコンテンツをマスクするために行われる。一例として、サーバ604は、猫のおもちゃを検索している、熱心な猫の愛好者に、犬のコンテンツを有するURLを返すように、検索エンジンをだますことができる。サーバ604は、正しいURLを使ってなりすましデータ614を編集するが、コンテンツは、猫に関連する情報に変更されている。しかし、実データ616は、犬に関連する情報を含む。したがって、クロラコンポーネント608は、実際には、犬に関する場合でも、URLが猫に関するものと思って、なりすましデータ614を検索する。クロラコンポーネント608は次いで、検索エンジンによってアクセス可能なデータ記憶コンポーネント610になりすましデータ614を格納する。その後、検索エンジンによる猫の検索も、犬の情報を含むURLを返す。この、犬/猫の例は害のないものと思われるであろうが、同じ技法が、たとえば広告、ポルノグラフィ、過激な文学、破壊活動集団、および他の主観的な攻撃的素材などをマスクするのに利用することができる。

#### 【0042】

図7に、本発明の態様による、ウェブクロラシステム702を伴うなりすまし防止処理700を示すブロック図を示してある。処理700は、ウェブクロラシステム702、クライアント制御コンポーネント704、ウェブサーバ706、および相互接続性を提供する通信システム708を含む。ウェブクロラシステム702は、クロラコンポーネント710、なりすましデータ714を有するデータ記憶コンポーネント712、および比較コンポーネント718を有する分散資源制御コンポーネント716からなる。クロラコンポーネント710は、図6に関して示し、説明したようなウェブサーバ706からなりすましデータ714を検索する。なりすましデータ714は次いで、データ記憶コンポーネント712に格納される。この時点で、一般的な検索エンジン(図7に示さず)は、なりすましデータ714にアクセスし、本当のコンテンツを知らない検索エンジンのユーザにデータ714を広める。しかし、本発明を利用することにより、なりすましデータ714を除去することができる。これは、サーバは、クロラに対してなりすましを行うが、通常、サーバのウェブページにアクセスするユーザに対してはなりすましを行わないという事実のおかげで遂行することができる。本発明では、クライアント制御コンポーネント704などの分散資源を利用するので、コンポーネント704は、サーバ706にユーザとしてアクセスし、サーバ706から実データを検索することができる。クライアント制御コンポーネント704は次いで、実データ(すなわち「クライアントデータ」)

10

20

30

40

50

および/または実データの表現を、分散資源制御コンポーネント716に転送することができる。分散資源制御コンポーネント716内部の比較コンポーネント718は次いで、格納されているなりすましデータ714を検索し、クライアントから受け取った実データと比較することができる。データが異なる場合、分散資源制御コンポーネント716は、データ記憶コンポーネント712にあるなりすましデータ714を上書きすることができる、その不正確さを排除する。こうすることにより、検索エンジンは、そうしないと利用可能にならなかった正確なデータにアクセスできるようになる。

#### 【0043】

上で示し説明した例示的なシステムを念頭において、図8~12のフロー図を参照すると、本発明によって実装することができる方法がよりよく理解されよう。説明を簡単にするために、この方法を一連のブロックとして示し、説明するが、本発明はブロックの順序に限定されないことを理解されたい。というのは、いくつかのブロックは、本発明によって、異なる順序で起こることもでき、かつ/または本明細書において示し説明する他のブロックと同時に起こることもできるからである。さらに、図示したすべてのブロックが、本発明による方法の実装に必要なわけではない。

#### 【0044】

本発明は、1つまたは複数のコンポーネント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という一般的な状況で説明することができる。概して、プログラムモジュールは、特定のタスクを実施または特定の抽象データタイプを実装するルーチン、プログラム、オブジェクト、データ構造などを含む。一般に、プログラムモジュールの機能は、様々な実施形態において所望される場合には、組み合わせることも分散することもできる。

#### 【0045】

図8を参照すると、本発明の態様による、クライアントベースのウェブクロウリングの方法800のフロー図を示してある。方法800は、802で始まり、804で、クライアントが、ウェブページを訪れたときに取得したウェブページ情報を記録する。本発明の単純な例において、情報は、訪れたウェブページのURLのみを含む。本発明のより複雑な例では、たとえば、URL、ウェブページのコンテンツのハッシュデータ、およびタイムスタンプなどを含むことができる。クライアントは次いで、806で、サーバにウェブページ情報を送る。本発明の一例では、クライアントが、他のクライアントにウェブページ情報を知らせることも可能である。やはり、単純な例では、情報は、URLのみを含むことができ、または複雑な例では、情報は、ウェブページについてのいくつかの異なるタイプのデータを含むことができる。本発明の一例では、クライアントは、ウェブページ情報から派生した付加情報を生成する。このデータは、たとえば、ウェブページがアクセス可能な期間、アクセスの容易さ(過負荷、接続再試行など)、および埋込みリンク状況などを含むことができる。さらに、ウェブページ情報は、いつ情報が送られるかを制御するようにスケジュールすることができる。スケジュールリングは、クライアントおよび/またはサーバによって開始することができる。情報を送るための基準は、時刻、時間の長さ、日付、集められたデータの量、および集められたデータのタイプ(たとえば、未知のデータまたは既知のデータ、発見されたなりすましデータ)などを含み得るが、それに限定されない。クライアントからサーバに送られる情報のボリュームを減らすために、本発明の一例では、クライアントが、ある特定のURLをサーバに知らせてあるかどうかに関する情報をローカルに保持し、まだ知らせていない場合のみ、サーバに情報を送ることができる。サーバは、ウェブページ情報を受け取ると、808で、情報を調べて、蓄積してあるデータと比較して新しいデータがあるかどうかを判定する。本発明の単純な例では、この判定は、サーバ上にすでに蓄積されたURL一覧と比較して、URLが新しいかどうか判定することを含む。未知の情報が見つかり、サーバは、810でサーバの蓄積データすなわち「既知の」データに未知の情報を追加し、812でフローが終わる。本発明の一例では、既知のデータは、サーバによる今後のウェブクロウリング、ダウンロード、および/または索引づけなどに利用されるURL一覧である。

10

20

30

40

50

## 【 0 0 4 6 】

図9に移ると、本発明の態様による、クライアントベースのウェブクロウリングの方法900の別のフロー図を示してある。方法900では、最初に、904で、クライアントが、訪れたウェブページのコンテンツに対するハッシュ値を計算する。クライアントが、そのウェブページを複数回訪れている場合、以前のハッシュ値が、そのウェブページに対して計算され、格納されている。クライアントは次いで、906で、計算したばかりの、すなわち「最新の」ハッシュ値を、ウェブページに対する以前のハッシュ値と比較する。比較を行っているクライアントは、見つかった違いの重要度を設定することができる。たとえば、ある場合には、ページ全体が変更されているが、別の場合には、ただ1つのコマがある文に追加されているだけである。クライアントは、変更の重要度を計算し、(a)この情報を使って、どの更新をサーバに送信するかという優先順位を決定するか、および/または(b)他のウェブページ情報とともに重要度の値をサーバに送信し、そうすることによって、サーバがページの再クロール/再索引づけに優先順位をつける際にこの情報を利用できるようにすることができる。変更重要度の例は、変更された文書の割合、変更の言語的/意味的重要度、および変更によって影響を受けるユーザ検索の割合の推定などのような項目を含むが、それに限定されない。重要度は、ページの人気の推定によって重みづけすることもできる。通常、クライアントは、908で、ウェブページ状況情報をローカルに格納し、必要な場合は、格納されたこの情報をアップデートする。クライアントは次いで、910で、ウェブページ状況情報をサーバに知らせる。本発明の一例では、クライアントが、他のクライアントにウェブページ状況情報を知らせることも可能である。サーバおよび/またはクライアントが通知を受ける方法は、URLのみ、URLと新しいハッシュ、ならびに/またはURLと新しいハッシュおよび古いハッシュなどを含み得るが、それに限定されない。サーバ(または他のクライアント)は次いで、クライアントのウェブページ状況情報が、URLの他に付加情報も含む場合、912で、その情報を、サーバ自体のウェブページ状況情報と比較する。サーバが、ウェブページのURLのみを状況変化として受信した場合、サーバは通常、サーバの以前のウェブページ状況情報と比較するための新しい状況情報を得るために、そのウェブページの再訪/クロールを開始する。サーバは次いで、914で、必要な場合は、サーバのウェブページ状況情報をアップデートし、916でフローが終わる。サーバがウェブページを再訪する必要をなくすために、クライアントは、付加情報を送信することができる。クライアントは、最後に訪れたときの、ページをキャッシュしたコピーをもっている場合、古いバージョンと新しいバージョンの間の違いとともに、古いハッシュ値、および新しいハッシュ値を送信することができる。サーバは最初に、クライアントの古いハッシュ値が、そのページの、サーバの現在のハッシュ値と一致するかを調べる。一致する場合、サーバは、それに従ってページのコンテンツをアップデートすることができる。

## 【 0 0 4 7 】

図10に、本発明の態様による、クライアントベースのウェブクロウリングの方法1000のさらに別のフロー図を示してある。方法1000は、1002で始まり、1004で、クライアントが、検索サーバ上で検索クエリを開始する。検索サーバは、1006で、検索クエリを分解し、クエリに回答して検索結果の一覧を編集する。検索サーバは次いで、1008で、埋込みウェブページリンク情報を有する検索結果ページを構成する。一般的な情報は、ウェブページのコンテンツの、サーバ側バージョンのハッシュ、および/または、コンテンツが、各ウェブページに対して新しくないと知られているかどうか(たとえば、別のクライアントが、ウェブページに対する新しいアップデートを検索サーバに知らせたが、検索サーバがそのページをまだ更新していない)を示すフラグを含み得るが、それに限定されない。したがって、「新しくないことが知られている」というフラグを有するウェブページは、検索サーバが、クライアントに、アップデート情報を送信するための検索を要求してほしくないウェブページである。検索サーバは次いで、1010で、検索を要求したクライアントに、埋込みリンクを有する検索結果ページを送信する。クライアントは、1012で、検索結果ページに列挙されているウェブページを訪れると、検

10

20

30

40

50

索サーバによって提供された、埋め込まれている鮮度フラグ（すなわち鮮度状況）を調べる。クライアントは、1014で、鮮度状況が「新鮮」であるとき、訪れたウェブページのコンテンツのハッシュを計算する。新鮮である状況は、検索サーバが、ウェブページの最近の、または最新のバージョンを所有していると思っ

10

てゐることを示す。したがって、クライアントは、ウェブページのコンテンツの新しいハッシュを計算し、1016で、このハッシュを、検索サーバによって提供された、埋め込まれているハッシュと比較する。クライアントは次いで、1018で、新しいハッシュと、検索サーバが提供したハッシュとの間に違いまたは差分が見つかったときは常に、検索サーバに通知する。検索サーバは次いで、1020で、通知を受信し、鮮度状況を「新しくないことが知られている」にアップデートし、再クロール用の一覧にもウェブページを追加し、1022でフローが終わ

【0048】

別の例では、上記の方法は、最新クライアント通知の時間、および最新クライアント通知からのウェブページのハッシュ値を含むがそれに限定されない、各検索結果を有する追加フィールドを検索サーバに送信させることによって、クライアントが、検索サーバにページ差異情報（この情報は後で、検索サーバが、ウェブページについてのサーバ側の情報を、クロールせずにアップデートするため、および/または検索サーバがウェブページを

20

いつ再クロールすべきかという優先順位づけを容易にするために使われる）を送信する方法を構成するように拡張される。クライアントが、検索サーバによって返されたページを訪れて、（a）「新鮮でないことが知られている」というフラグが偽であるか、または（b）「新鮮でないことが知られている」が真であり、かつ最新のクライアント通知からのハッシュ値が、このクライアントがウェブページに対して計算したハッシュ値と異なる場合、クライアントは、検索サーバに通知を行う。ページ変更の周期を認識することによって、ページがAからB、C、Aへと繰り返し変更される場合、その変更を認識し、このページについてのクライアントによるアップデートを制限できるようにすることも可能である。

【0049】

検索結果ページを使うことによって、クライアントとサーバの間のトラフィックは、本発明によって示したように、クライアントベースのウェブクロールにおいて徹底して削減される。さらに、クライアントの秘密は、サーバが検索結果ページ中で提供したウェブページのみがクライアントによってアップデートされることにより保たれる。このようにして、クライアントが、限定アクセスを有するウェブページを訪れた場合、この情報は、不注意によって検索サーバに送信されない。この方法の利点は、新しいウェブページが秘密でない場合でも、検索サーバが、クライアントを用いて、検索において利用するために既知のウェブページを拡張できないという点で、欠点でもある。

【0050】

本発明の別の例では、方法（図示せず）は、検索サーバのウェブページ情報を利用するだけでなく、他の検索サーバのウェブページ情報も利用する。したがって、別の検索サーバによってクライアントに返された新しいウェブページは、新しいウェブページが存在することを検索サーバに通知するのに使うことができる。この方法もやはり、クライアントの秘密を保つ。というのは、この方法は、検索サーバが列挙していない、公に利用可能なウェブページのみを検索サーバに通知するからである。こうすることにより、検索サーバは、クライアントの信用を損なわずに未知のウェブページを追加できるようになる。通知は、ウェブページのURL、ウェブページのコンテンツのハッシュ、ウェブページにアクセスしたときのタイムスタンプ、およびウェブページに対する以前のハッシュと比較した、新しいハッシュの差分などを含み得るが、それに限定されない。

【0051】

10

20

30

40

50

図11に移ると、本発明の態様による、クライアントベースのウェブクローリングの方法1100のさらに別のフロー図を示してある。方法1100は、1102で始まり、1104で、検索サーバが、弱（損失を伴う）インジケータ関数の集合を生成する。こうした関数を生成する方法は、後で説明する。検索サーバは、1106で、クライアントベースのウェブクロウラを備えるクライアントに、ランダムに選択された弱インジケータ関数を送信する。クライアントは次いで、1108で、ランダムに選択された弱インジケータ関数によって既知でないこと示されるウェブページ用のウェブページのデータを生成する。概して、既知でないウェブページのみが、弱インジケータ関数によって正確に表される。「既知の」ウェブページは、実際に既知であっても、既知でなくてもよい。クライアントは次いで、1110で、未知のウェブページのデータをサーバに送信する。サーバは次いで、1112で、このデータを使って、ウェブページに関するサーバ側の情報をアップデートし、1114でフローが終わる。

#### 【0052】

図12を参照すると、本発明の態様による、クライアントベースのウェブクローリングのための弱インジケータ関数の適切な集合を生成する方法1200のフロー図を示してある。方法1200は、1202で始まり、1204で、検索サーバ上にあるウェブページ情報を表す辞書を、重なりのない部分辞書にランダムに分割する。一般に、部分辞書は、それぞれがウェブページ情報のグループ化における共通の特色を表すように選ばれる。弱（損失を伴う）インジケータ関数は次いで、1206で、各部分辞書が、ある特定の部分辞書にあるウェブページ情報を表すようにランダムに選ばれる。関数は次いで、1208で、少なくとも1つの部分辞書の弱関数が1に等しい場合、かつその場合に限り、 $I(x) = 1$ となるように作成され、1210でフローが終わる。このようにして、弱インジケータ関数の「適切な集合」が生成される。したがって、たとえば、サーバによって知られているどのURLに対しても、インジケータ関数は、URLが既知であると正しく判定する。サーバによって知られていないURLに対して、インジケータ関数は、既知であると誤ってラベルづけする可能性があるが、その場合、クライアントは、何もしないか、または未知であると正しくラベルづけする可能性があり、この場合、クライアントはサーバに知らせることができる。定義により、サーバにとって未知であるウェブサイトがクライアントによって訪問されるときはいつでも、クライアントのインジケータ関数がそのサイトを新しいものであると認識する確率がゼロでないことを、弱インジケータ関数の適切な集合が保証する。

#### 【0053】

本発明の様々な態様を実装するさらなる状況を提供するために、図13および以下の説明では、本発明の様々な態様を実装することができる適切な計算機環境1300の、簡潔で全般的な説明を提供することを意図している。これまでは、ローカルコンピュータおよび/またはリモートコンピュータを実行するコンピュータプログラムのコンピュータ実行可能命令という一般的な状況において本発明を説明したが、本発明は他のプログラムモジュールとの組合せにおいても実装できることが当業者には理解されよう。概して、プログラムモジュールは、特定のタスクを実施し、かつ/または特定の抽象データタイプを実装するルーチン、プログラム、コンポーネント、データ構造などを含む。さらに、発明性のある本方法は、他のコンピュータシステム構成とともに実施できることが当業者には理解されよう。他のコンピュータシステム構成には、シングルプロセッサコンピュータシステムまたはマルチプロセッサコンピュータシステム、ミニコンピュータ、メインフレームコンピュータ、ならびにパーソナルコンピュータ、ハンドヘルド計算装置、マイクロプロセッサベースの家電製品および/またはプログラム可能な家電製品などがあり、それぞれが1つまたは複数の関連する装置と動作可能に通信することができる。図示した本発明の態様は、通信ネットワークを介してリンクされるリモート処理装置によって特定のタスクが実施される分散型計算機環境でも実施することができる。ただし、すべてではなくともいくつかの本発明の態様は、独立型のコンピュータにおいても実施することができる。分散型計算機環境では、プログラムモジュールは、ローカルメモリ記憶装置および/またはリ

10

20

30

40

50

モートメモリ記憶装置内に配置することができる。

【0054】

本アプリケーションで使用する「コンポーネント」という用語は、ハードウェア、ハードウェアおよびソフトウェアの組合せ、ソフトウェア、または実行中のソフトウェアのいずれかであるコンピュータ関連エンティティを指すことを意図している。たとえば、コンポーネントは、プロセッサで実行中の処理、プロセッサ、オブジェクト、実行ファイル、実行のスレッド、プログラム、およびコンピュータでよいが、それに限定されない。実例として、サーバ上で実行中のアプリケーションおよび/またはそのサーバがコンポーネントとなり得る。さらに、コンポーネントは、1つまたは複数の下位コンポーネントを含むことができる。

10

【0055】

図13を参照すると、本発明の様々な態様を実装する例示的なシステム環境1300は、従来のコンピュータ1302を含み、このコンピュータは、処理装置1304、システムメモリ1306、およびシステムメモリを含む様々なシステムコンポーネントを処理装置1304に結合するシステムバス1308を含む。処理装置1304は、市販されているどのプロセッサでも、固有のどのプロセッサでもよい。さらに、この処理装置は、並列に接続することができるような複数のプロセッサから形成されるマルチプロセッサとして実装することができる。

【0056】

システムバス1308は、従来の様々なバスアーキテクチャ、たとえばいくつか例を挙げると、PCI、VESA、マイクロチャンネル、ISA、およびEISAのどれをも使用するメモリバスまたはメモリコントローラ、周辺バス、およびローカルバスなどいくつかのタイプのバス構造のいずれでもよい。システムメモリ1306は、ROM（読み出し専用メモリ）1310およびRAM（ランダムアクセスメモリ）1312を含む。BIOS（基本入出力システム）1314は、たとえば起動中にコンピュータ1302内部の要素間の情報の転送を助ける基本ルーチンを含み、ROM1310に格納される。

20

【0057】

コンピュータ1302は、たとえば、ハードディスクドライブ1316、たとえば取外し可能ディスク1320からの読み出しまたはそこへの書込みを行うための磁気ディスクドライブ1318、および、たとえばCD-ROMディスク1324または他の光学媒体からの読み出しまたはそこへの書込みを行う光ディスクドライブ1322を含むこともできる。ハードディスクドライブ1316、磁気ディスクドライブ1318、および光ディスクドライブ1322は、それぞれハードディスクドライブインターフェイス1326、磁気ディスクドライブインターフェイス1328、および光ドライブインターフェイス1330によって、システムバス1308に接続される。ドライブ1316~1322およびそれに関連するコンピュータ可読媒体は、データ、データ構造、コンピュータ実行可能命令などを含む不揮発性記憶をコンピュータ1302に提供する。上記のコンピュータ可読媒体の説明では、ハードディスク、取外し可能な磁気ディスク、およびCDに言及したが、コンピュータ可読な他のタイプの媒体、たとえば磁気カセット、フラッシュメモリカード、デジタル映像ディスク、ベルヌーイカートリッジなども、例示的な動作環境1300において使うことができ、さらに、このような媒体も、本発明の方法を実施するコンピュータ実行可能命令を含むことができることが当業者には理解されよう。

30

40

【0058】

オペレーティングシステム1332、1つまたは複数のアプリケーションプログラム1334、他のプログラムモジュール1336、およびプログラムデータ1338などいくつかのプログラムモジュールを、ドライブ1316~1322およびRAM1312に格納することができる。オペレーティングシステム1332は、適切などのオペレーティングシステムでも、オペレーティングシステムの組合せでもよい。一例として、アプリケーションプログラム1334およびプログラムモジュール1336が、本発明の態様によるクライアントベースのウェブクロールリングを容易にすることを含むことができる。

50

## 【0059】

ユーザは、キーボード1340およびポインティングデバイス（たとえばマウス1342）など1つまたは複数のユーザ入力装置を介して、コマンドおよび情報をコンピュータ1302に入力することができる。他の入力装置（図示せず）には、マイクロホン、ジョイスティック、ゲーム用パッド、衛星パラボラアンテナ、無線リモコン、スキャナなどがあり得る。こうしたおよび他の入力装置はしばしば、システムバス1308に結合されるシリアルポートインターフェイス1344を介して処理装置1304に接続されるが、他のインターフェイス、たとえばパラレルポート、ゲームポート、またはUSB（ユニバーサルシリアルバス）によって接続することもできる。モニタ1346または他のタイプの表示装置も、ビデオアダプタ1348などのインターフェイスを介してシステムバス1308に接続される。モニタ1346に加えて、コンピュータ1302は、他の周辺出力装置（図示せず）、たとえばスピーカ、プリンタなども含むことができる。

10

## 【0060】

コンピュータ1302は、1つまたは複数のリモートコンピュータ1360への論理接続を使用してネットワーク接続された環境において動作できることを理解されたい。リモートコンピュータ1360は、ワークステーション、サーバコンピュータ、ルータ、ピア装置、または他の共通ネットワークノードでよく、通常、コンピュータ1302に関連して説明した要素の多くまたはすべてを含むが、簡潔にするために、メモリ記憶装置1362のみを図13に示した。図13に示した論理接続は、LAN（ローカルエリアネットワーク）1364およびWAN（ワイドエリアネットワーク）1366を含むことができる。このようなネットワーク環境は、オフィス、企業規模のコンピュータネットワーク、イントラネット、およびインターネットにおいてよく見られる。

20

## 【0061】

LANネットワーク環境において使われる場合、たとえば、コンピュータ1302は、ネットワークインターフェイスまたはアダプタ1368を介してローカルネットワーク1364に接続される。WANネットワーク環境において使われる場合、コンピュータ1302は通常、モデム（たとえば、電話、DSL、ケーブルなど）1370を含み、またはLAN上の通信サーバに接続され、あるいは、たとえばインターネットなどのWAN1366を介した通信を確立する他の手段を有する。モデム1370は、コンピュータ1302に対して内蔵型でも外付け型でもよく、シリアルポートインターフェイス1344を介してシステムバス1308に接続される。ネットワーク接続された環境では、プログラムモジュール（アプリケーションプログラム1334など）および/またはプログラムデータ1338は、リモートメモリ記憶装置1362に格納することができる。図示したネットワーク接続は例示的なものであり、本発明の態様を実施する際に、コンピュータ1302と1360の間の通信リンクを確立する他の手段（たとえば、有線または無線）も使うことができることが理解されよう。

30

## 【0062】

他の指示がない限り、コンピュータプログラミングの当業者による実施に従って、コンピュータ、たとえばコンピュータ1302またはリモートコンピュータ1360によって実施される作用および象徴的に表した動作を参照して本発明を説明した。このような作用および動作は、コンピュータに実行されるものとして何度か言及した。こうした作用および象徴的に表した動作は、処理装置1304による、データビットを表す電気信号の処理を含み、その結果、電気信号表示の変換または減少、およびメモリシステム（システムメモリ1306、ハードドライブ1316、フロッピー（登録商標）ディスク1320、CD-ROM1324、およびリモートメモリ1362など）内のメモリ位置でのデータビットの保持を引き起こし、そうすることによって、コンピュータシステムの動作、ならびに他の信号処理を再構成し、あるいは変更することが理解されよう。このようなデータビットが保持されるメモリ位置は、データビットに対応する特定の電気属性、磁気属性、または光学属性を有する物理的な位置である。

40

## 【0063】

50

図14は、本発明と相互作用する一例である計算機環境1400の別のブロック図である。システム1400はさらに、1つまたは複数のクライアント(群)1402を含むシステムを示す。クライアント(群)1402は、ハードウェアおよび/またはソフトウェア(たとえば、スレッド、処理、計算装置)でよい。システム1400は、1つまたは複数のサーバ(群)1404も含む。サーバ(群)1404は、ハードウェアおよび/またはソフトウェア(たとえば、スレッド、処理、計算装置)でよい。サーバ1404は、たとえば、本発明を利用して変換を実施するためのスレッドを収容することができる。クライアント1402とサーバ1404の間の可能な1つの通信は、2つ以上のコンピュータ処理の間で伝送されるように適合されたデータパケットの形をとることができる。システム1400は、クライアント(群)1402とサーバ(群)1404の間の通信を容易にするのに利用できる通信フレームワーク1408を含む。クライアント(群)1402は、クライアント(群)1402にローカルな情報を格納するのに利用できる、1つまたは複数のクライアントデータストア(群)1410に動作可能に接続される。同様に、サーバ(群)1404は、サーバ1404にローカルな情報を格納するのに利用できる、1つまたは複数のサーバデータストア(群)1406に動作可能に接続される。

10

**【0064】**

本発明の一例では、ウェブクロールを容易にする、2つ以上のコンピュータコンポーネントの間で伝送されるデータパケットは、少なくとも部分的に、ウェブクロール用の分散型システムを少なくとも部分的に使用するウェブクロールに関する情報からなる。

20

**【0065】**

本発明の別の例では、ウェブクロールを容易にするシステムの、コンピュータ実行可能なコンポーネントを格納するコンピュータ可読媒体は、少なくとも部分的には、ウェブクロール用の分散システムによって編集される、ウェブページに関連する情報を少なくとも部分的に判定するウェブクロールシステムからなる。

**【0066】**

本発明のシステムおよび/または方法は、コンピュータコンポーネント、および非コンピュータ関連コンポーネントを同様に容易にするウェブクロールシステムにおいて利用できることを理解されたい。さらに、本発明のシステムおよび/または方法は、有線および/または無線などでよい、コンピュータ、サーバ、および/またはハンドヘルド電子装置などを含むがそれに限定されない広範囲の電子関連技術において利用可能であることが当業者には理解されよう。

30

**【0067】**

本発明は、サーバ-クライアントベースのクロールシステムだけでなく、ピアツーピアのクロールシステムにも利用できることも当業者には理解されよう。クライアントは、一般に「サーバ」の動作に関連づけられたタスクを実施することができ、したがって、本発明のいくつかの例において、サーバに関連付けられたいくつかの特性をクライアントに転送することも可能である。本発明の一事例として、クライアントは、他のクライアントに対して「部分クロール」を実施して、サーバに送信するための情報を確認し、かつ/または検索する。この例は、たとえば、特定のクライアントとサーバの間のボトルネックを有するネットワークにおいて有益であり得る。データは、サーバへの最高のアクセス権を有するクライアントに転送することができる。本発明の他の例では、クライアントは、イントラネットシステムにおいて部分クロールを開始することによってサーバの動作を示すことができ、したがって、イントラネット上に存在する唯一の、および/または大幅に削減された数のクライアントから、サーバに情報を報告する。このようにして、検索サーバは、クライアントにおいて様々な部分クロールを開始して、サーバのクロール用資源を拡張することができる。

40

**【0068】**

上記の説明内容は、本発明のいくつかの例を含む。当然ながら、本発明を説明するためのコンポーネントまたは方法のあらゆる組合せを説明することはできないが、本発明のさ

50

らに多くの組合せおよび入替えが可能であることが当業者には理解できよう。したがって、本発明は、添付の特許請求の範囲の精神および範囲内であるこのようなすべての変更形態、修正形態、および変形形態を包含することを意図したものである。さらに、詳細な説明または特許請求の範囲において「含む」という用語が使われている限りでは、「備える」は、使用される場合、請求項においては接続語として解釈されるが、「含む」のような用語は、「備える」という用語と同様に包括的であることを意図している。

【図面の簡単な説明】

【0069】

【図1】本発明の態様によるデータ分析システムを示すブロック図である。

【図2】本発明の態様によるデータ分析システムを示す別のブロック図である。

10

【図3】本発明の態様によるデータ分析システムを示すさらに別のブロック図である。

【図4】本発明の態様によるデータ分析システムを示すさらに別のブロック図である。

【図5】本発明の態様による、ページ検索結果を使用するデータ分析システムを示す図である。

【図6】本発明の態様による、ウェブクローラシステムを伴うなりすまし処理を示すブロック図である。

【図7】本発明の態様による、ウェブクローラシステムを伴うなりすまし防止処理を示すブロック図である。

【図8】本発明の態様による、クライアントベースのウェブクローリングの方法を示すフロー図である。

20

【図9】本発明の態様による、クライアントベースのウェブクローリングの方法を示す別のフロー図である。

【図10】本発明の態様による、クライアントベースのウェブクローリングの方法を示すさらに別のフロー図である。

【図11】本発明の態様による、クライアントベースのウェブクローリングの方法を示すさらに別のフロー図である。

【図12】本発明の態様による、クライアントベースのウェブクローリングのための弱インジケータ関数の適切な集合を生成する方法を示すフロー図である。

【図13】本発明が機能することができる一例の動作環境を示す図である。

【図14】本発明が機能することができる別の例の動作環境を示す図である。

30

【符号の説明】

【0070】

100 データ分析システム

102 ~ 106 クライアント

110 検索サーバ

112 ウェブページサーバ

200 データ分析システム

202 クライアント

204 サーバ

300 データ分析システム

40

400 データ分析システム

500 データ分析システム

502 クライアント

504 サーバ

508 サーバへ送信

510 サーバから受信

600 なりすまし処理

700 なりすまし防止処理

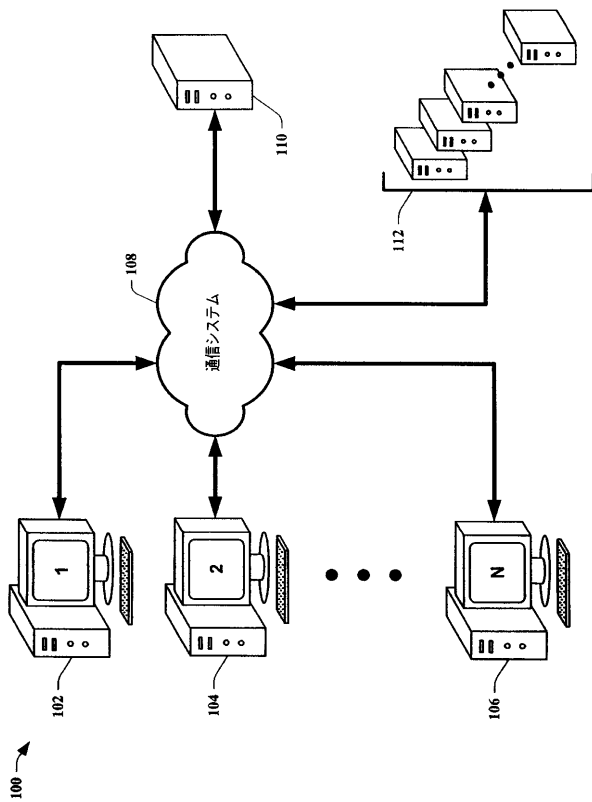
706 ウェブサーバ

800 クライアントベースのウェブクローリングの方法

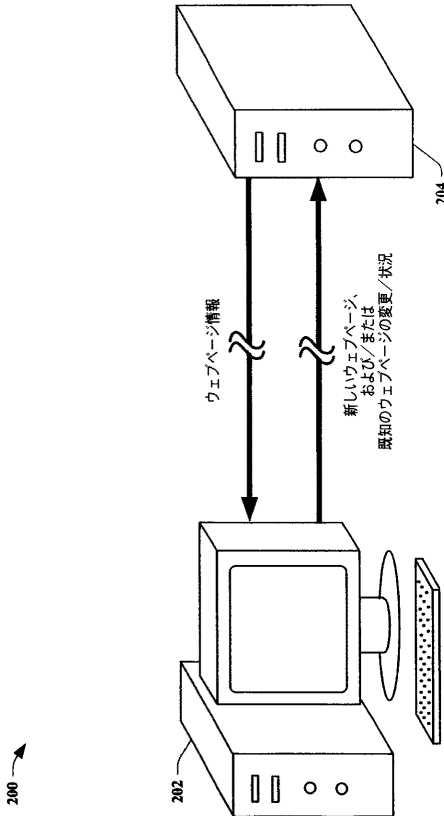
50

- 9 0 0 クライアントベースのウェブクロージングの方法
- 1 0 0 0 クライアントベースのウェブクロージングの方法
- 1 1 0 0 クライアントベースのウェブクロージングの方法
- 1 2 0 0 弱インジケータ関数の適切な集合を生成する方法
- 1 3 0 0 例示的なシステム環境
- 1 3 0 2 従来のコンピュータ
- 1 3 1 2 R A M
- 1 3 1 6 ハードディスクドライブ
- 1 3 1 8 磁気ディスクドライブ
- 1 3 2 0 取外し可能ディスク
- 1 3 2 2 光ディスクドライブ
- 1 3 2 4 C D - R O M ディスク
- 1 3 4 0 キーボード
- 1 3 4 2 マウス
- 1 3 6 2 メモリ記憶装置
- 1 4 0 0 計算機環境

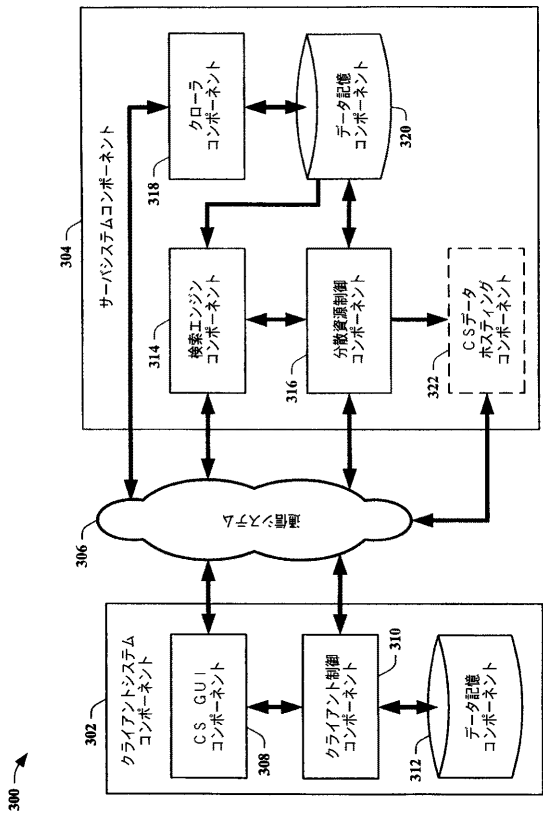
【 図 1 】



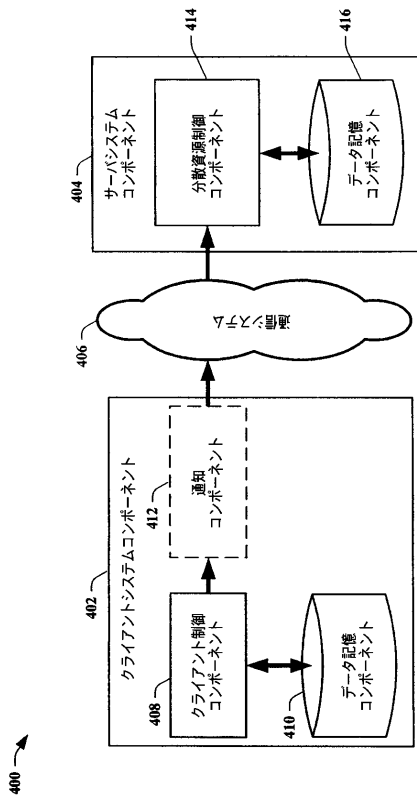
【 図 2 】



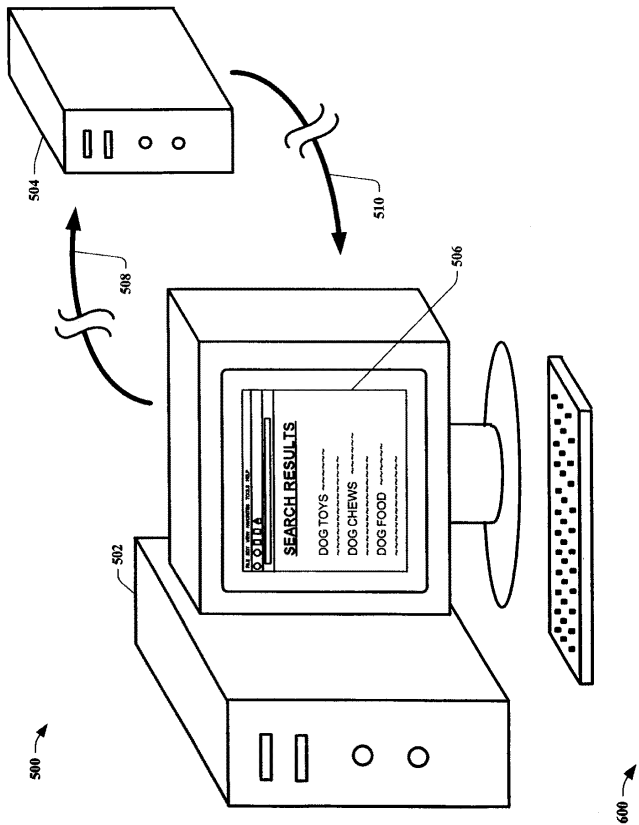
【 図 3 】



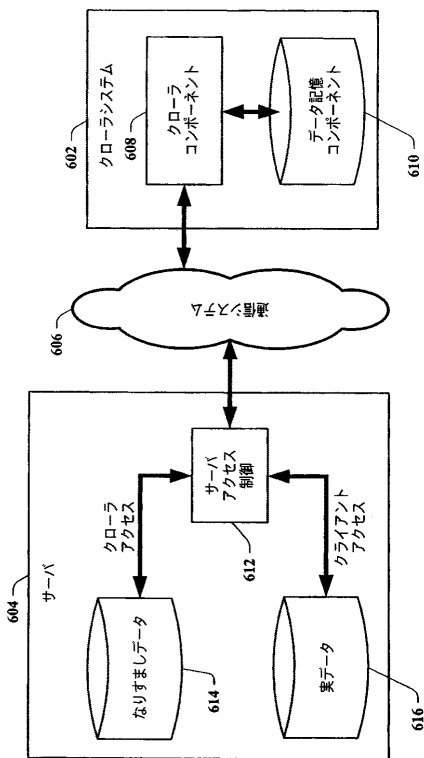
【 図 4 】



【 図 5 】

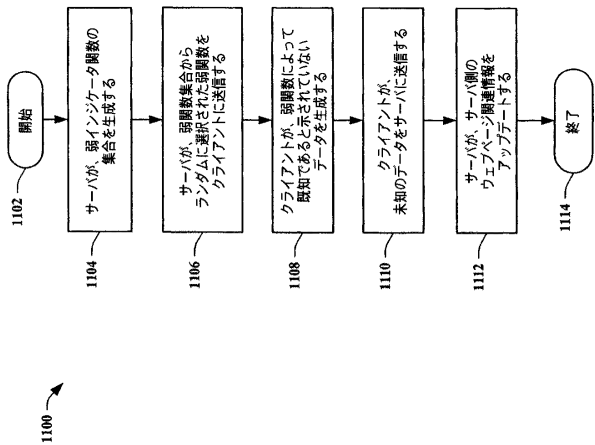


【 図 6 】

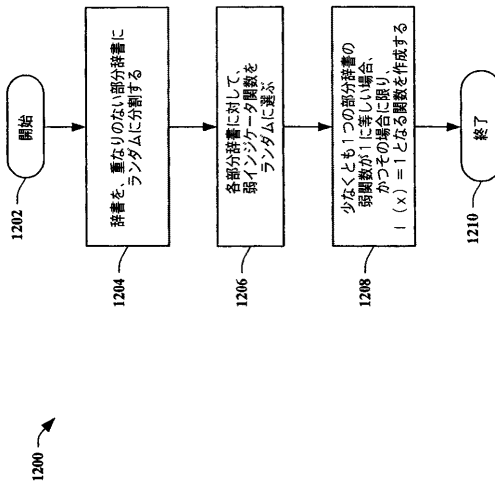




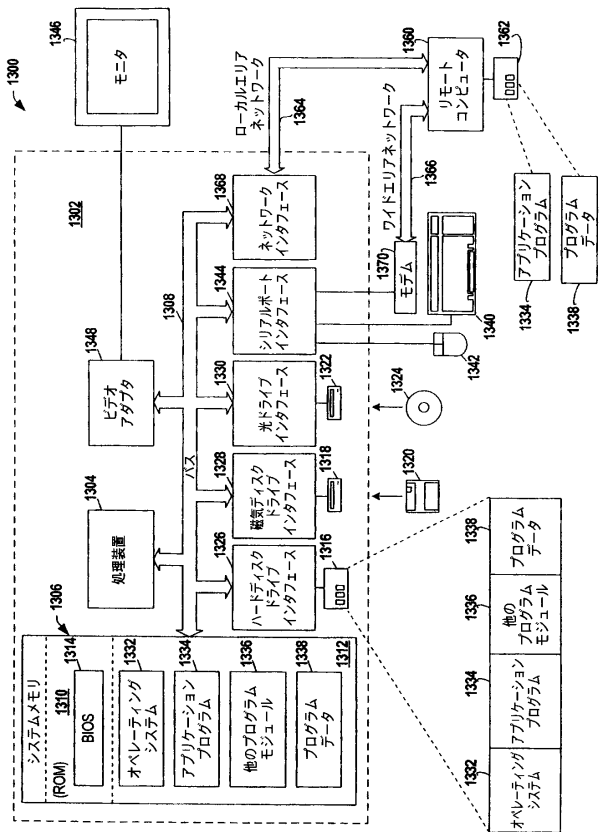
【図 1 1】



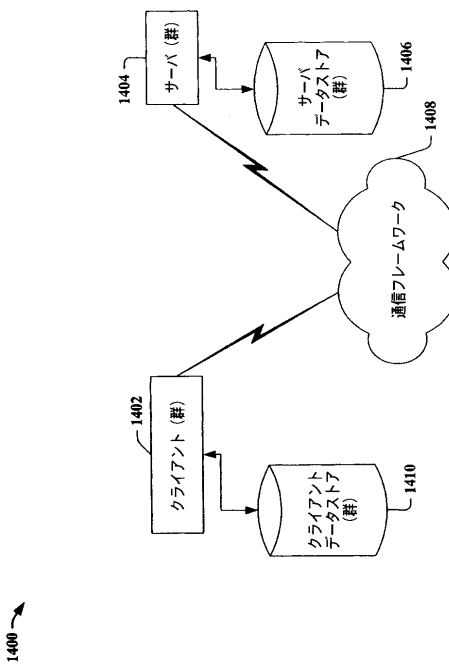
【図 1 2】



【図 1 3】



【図 1 4】



---

フロントページの続き

(72)発明者 クリストファー エー・ミーク  
アメリカ合衆国 98033 ワシントン州 カークランド ノースイースト 71 ストリート  
12935

審査官 田上 隆一

(56)参考文献 特開2002-312284(JP,A)  
特開2002-140257(JP,A)  
岡崎 浩基, 訪問者を確実に増やす 誰でもできるアクセス解析読み取り術 眠れるお宝を活用  
しろ!, 月刊CYBIZ SOHO computing, 日本, 株式会社サイビズ, 2003  
年 5月 1日, 第8巻 第5号

(58)調査した分野(Int.Cl., DB名)  
G06F 13/00