

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3896014号
(P3896014)

(45) 発行日 平成19年3月22日(2007.3.22)

(24) 登録日 平成18年12月22日(2006.12.22)

(51) Int. Cl.	F I
G06F 17/30 (2006.01)	G06F 17/30 340A
G06F 13/00 (2006.01)	G06F 17/30 180Z
	G06F 17/30 380D
	G06F 13/00 560A

請求項の数 6 (全 24 頁)

(21) 出願番号	特願2002-81642 (P2002-81642)	(73) 特許権者	000003078
(22) 出願日	平成14年3月22日 (2002.3.22)		株式会社東芝
(65) 公開番号	特開2003-281173 (P2003-281173A)		東京都港区芝浦一丁目1番1号
(43) 公開日	平成15年10月3日 (2003.10.3)	(74) 代理人	100058479
審査請求日	平成15年4月17日 (2003.4.17)		弁理士 鈴江 武彦
		(74) 代理人	100084618
			弁理士 村松 貞男
		(74) 代理人	100092196
			弁理士 橋本 良郎
		(74) 代理人	100091351
			弁理士 河野 哲
		(74) 代理人	100088683
			弁理士 中村 誠
		(74) 代理人	100070437
			弁理士 河井 将次

最終頁に続く

(54) 【発明の名称】 情報収集システム、情報収集方法及びコンピュータに情報収集を実行させるプログラム

(57) 【特許請求の範囲】

【請求項1】

ユーザの要求を満足する情報を収集して提示する情報収集システムにおいて、
それぞれ複数のユーザをメンバーとする複数のコミュニティを管理するコミュニティ管理手段と、

各コミュニティに属するメンバーがメッセージの送受信を行うためのメッセージ送受信手段と、

前記複数のコミュニティの各々で共有されている情報をユーザが閲覧するためのコミュニティ情報提示手段と、

各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集要求として、
収集の起点とする情報と、情報が含むべき語句の条件とを記述した収集要求を編集するための収集要求編集手段と、

各複数のコミュニティにおいて編集された複数の収集要求のいずれかを満足する情報を、
前記各々の収集要求に記述された収集の起点である情報からハイパーリンクを辿って前記語句の条件を満たす情報を探索することにより、情報ネットワーク上の複数の情報源から収集する情報収集手段と、

前記収集した情報に基づいて前記複数の収集要求の各々に対応する収集結果を各々生成する収集結果生成手段と、

各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集結果を編集するための収集結果編集手段と、を具備し、

10

20

前記コミュニティ情報提示手段は、複数のコミュニティで各々作成された複数の収集結果を、当該コミュニティまたは当該コミュニティ内で送受信されるメッセージと関連付けて、当該コミュニティのメンバーおよび非メンバーのユーザに提示すると共に、当該コミュニティにおける当該収集結果を構成する前記情報収集手段で収集された情報と、他のコミュニティにおける収集結果を構成する前記情報収集手段で収集された情報とが、重複する場合には、当該コミュニティにおける当該収集結果と、当該他のコミュニティにおける収集結果とを、関連付けて提示することを特徴とする情報収集システム。

【請求項2】

請求項1に記載の情報収集システムにおいて、コミュニティのメンバーが前記メッセージ送受信手段を用いて送受信するメッセージから、前記情報収集の起点とし得る情報と、当該情報に関わるコメント文とを抽出して、これらに基づき、当該コミュニティの収集要求及び当該コミュニティの収集結果の少なくとも一方を自動的に更新することを特徴とする情報収集システム。

10

【請求項3】

請求項1又は請求項2に記載の情報収集システムにおいて、コミュニティのメンバーが前記収集結果編集手段を用いて行った収集結果の編集内容に基づき、当該収集結果に対応する収集要求を更新することを特徴とする情報収集システム。

【請求項4】

請求項1から請求項3のいずれか1項に記載の情報収集システムにおいて、ユーザが入力する検索条件を満足する情報を、前記情報収集手段で収集した情報の中から検索する収集情報検索手段をさらに具備し、当該収集情報検索手段は、検索された情報と、コミュニティで作成した収集結果のうち前記検索された情報を含む収集結果とを、関連付けて提示することを特徴とする情報収集システム。

20

【請求項5】

ユーザの要求を満足する情報を収集して提示する情報収集方法において、コンピュータが、各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集要求として、収集の起点とする情報と、情報が含むべき語句の条件とを記述した収集要求を編集し、

コンピュータが、各コミュニティにおいて編集された複数の収集要求のいずれかを満足する情報を、前記各々の収集要求に記述された収集の起点である情報からハイパーリンクを辿って前記語句の条件を満たす情報を探索することにより、情報ネットワーク上の複数の情報源から収集し、

30

コンピュータが、前記収集した情報に基づいて前記複数の収集要求の各々に対応する収集結果を各々生成し、

コンピュータが、各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集結果を編集し、

コンピュータが、複数のコミュニティが各々作成された複数の収集結果を、当該コミュニティまたは当該コミュニティ内で送受信されるメッセージと関連付けて、前記複数のコミュニティの各々で共有されている情報を、当該コミュニティのメンバーおよび非メンバーのユーザに提示すると共に、当該コミュニティにおける当該収集結果を構成する前記情報収集手段で収集された情報と、他のコミュニティにおける収集結果を構成する前記情報収集手段で収集された情報とが、重複する場合には、当該コミュニティにおける当該収集結果と、当該他のコミュニティにおける収集結果とを、関連付けて提示することを特徴とする情報収集方法。

40

【請求項6】

コンピュータにユーザの要求を満足する情報を収集して提示する情報収集を実行させるプログラムにおいて、

コンピュータに、各コミュニティに属するメンバーによって共同で編集された当該コミュニティにおける収集要求として、収集の起点とする情報と、情報が含むべき語句の条件とを記述した収集要求を入力させ、

50

コンピュータに、各コミュニティにおいて編集された複数の収集要求のいずれかを満足する情報を、前記各々の収集要求に記述された収集の起点である情報からハイパーリンクを辿って前記語句の条件を満たす情報を探索することにより、情報ネットワーク上の複数の情報源から収集させ、

コンピュータに、前記収集した情報に基づいて前記複数の収集要求の各々に対応する収集結果を各々生成させ、

コンピュータに、各コミュニティに属するメンバーが共同で編集された当該コミュニティにおける収集結果を入力させ、

コンピュータに、複数のコミュニティで各々作成された複数の収集結果を、当該コミュニティまたは当該コミュニティ内で送受信されるメッセージと関連付けて、前記複数のコミュニティの各々で共有されている情報を、当該コミュニティのメンバーおよび非メンバーのユーザに提示させると共に、当該コミュニティにおける当該収集結果を構成する前記情報収集手段で収集された情報と、他のコミュニティにおける収集結果を構成する前記情報収集手段で収集された情報とが、重複する場合には、当該コミュニティにおける当該収集結果と、当該他のコミュニティにおける収集結果とを、関連付けて提示させることを特徴とするプログラム。

10

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、インターネットやイントラネットなどの情報ネットワークに分散して配置された複数の情報源から、ユーザの要求を満足する情報を収集する情報収集システムに関する。

20

【0002】

【従来の技術】

大規模情報ネットワーク技術の普及により、誰もが大量の情報を自由に利用できるようになっている。インターネットやイントラネットなどの情報ネットワーク上には、大量の情報がハイパーテキスト形式の文書（ウェブページ）として公開されており、その数は数十億とも言われている。これらの情報を利用する方法として、ブラウザと呼ばれる情報閲覧ソフトウェアを用い、関心のある項目（ハイパーリンク）を選択していく（ブラウジングする）方法が用いられるのが一般的である。また、大量の情報の中から、キーワード等で指定した条件を満足する情報を検索するための検索サービスサイトや、情報を利用しやすい形に分類して提供するディレクトリサイトが、各種運用されている。ユーザは、所望の情報を得るためには、まず、検索サービスサイトやディレクトリサイトを用いて自分の関心に合致しそうな文書を求めた後に、その文書の内容や、その文書にリンクされた他の文書の内容を、ブラウジングすることによって調べるといって、一連の作業を繰返し行う。また、頻繁に利用する情報や特に重要な情報については、ブラウザの付属機能であるブックマークと呼ばれる手段を用いてその情報の位置（URL）を記憶したり、有用な情報の位置をリストアップした文書（リンク集）を作成して利用することが行われている。

30

【0003】

しかし、大量の情報の中から、検索やブラウジングによって必要な情報を集める作業は時間と労力を要する。また、検索サービスサイトやディレクトリサイトでは、最新の情報や、専門性の高い情報が十分に提供されていないという問題もある。これらの問題を解決する技術の一つに、自動クロウリング技術が知られている。これは、ハイパーテキストのハイパーリンクを再起的に辿る（すなわち、クロウリングする）ソフトウェア（すなわち、クローラ）を用いて、大量の文書情報を自動的に走査し、ユーザが指定した条件を満足する文書を収集する方法である。ユーザがクローラに与えることのできる収集条件には、収集する文書の個数・容量の制限や、収集を開始する起点の文書、起点の文書から辿るリンクの段数の上限、収集する範囲（ウェブサーバのドメインなど）、文書の更新日時の条件、などがある。また、文書の内容に関する条件としては、キーワード・フレーズ等が対象文書中で出現する頻度や、例示した文書と対象文書との類似度、ユーザの興味・関心の記

40

50

述（プロフィール）と対象文書との類似度、などについての条件がある。さらには、対象文書の重要度を、アクセス数やハイパーリンクの構造に基づいて計算し、重要度の大きい文書を優先的に収集する方法なども提案されている。自動クロール技術に関する公知文献には、“Focused Crawling: A New Approach for Topic-Specific Resource Discovery”, Soumen Chakrabarti他, The Eighth International World Wide Web Conference, 1999（以下、「文献1」と称する）や、特開平10-260978号公報「情報収集方法及び装置」（以下、「文献2」と称する）などがある。

【0004】

一方、複数のユーザが互いに情報を交換するための手段としては、電子メールおよびメーリングリスト、電子掲示板、チャットなどの手段が、広く普及している。メーリングリストは、複数のユーザの電子メールアドレスをまとめて、その全員に一括してメッセージを送信できるようにした手段である。また、電子掲示板は、ネットワーク上に情報共有のためのスペースを設けて、複数の登録ユーザあるいは匿名ユーザが自由にメッセージを記入できるようにした手段である。チャットは、電子掲示板と同様に情報共有スペースを設けて、テキストのメッセージをリアルタイムに送受信できるようにした手段である。メーリングリストや電子掲示板、チャット等のように、比較的多数のユーザによる（一対一のみでない）メッセージの交換を目的としたコミュニケーション手段では、参加メンバーの大部分が共通に関心を持つ話題に関するメッセージがやり取りされることが多い。このように、共通の目的や話題を持って電子的なメッセージを交換するユーザの集団を、本明細書においては、以下、「コミュニティ」と称する。

【0005】

コミュニティのメンバーの一人が有用な情報を得た場合、上述のコミュニケーション手段を用いて他のメンバーに通知することによって、メンバー間で情報を共有するということが日常的に行われている。このようにして交換される情報のうち、とくに有用な情報については、メンバーの有志が自発的に、有用な情報を手作業でリストアップし、他のメンバーが利用しやすいようにリンク集などの形に整理し、定期的に保守するということが行われる場合もある。コミュニティのメンバーが関心を持つ話題は、コミュニティの趣旨を逸脱しない範囲内にある場合が多いが、多少は動的に変遷する。コミュニティのメンバーがどのような話題に関心を持っているかを自動的に調べる技術については、特開2000-293526号公報「嗜好情報収集システム」（以下、「文献3」と称する）や、特開2001-92755号公報「プロフィール作成方法及びシステム」（以下、「文献4」と称する）などの公知文献がある。

【0006】

【発明が解決しようとする課題】

自動クロールは、収集に要する時間とネットワーク資源の消費が大きいわりに収集の効率が良くないという問題がある。インターネットからのクロールによる収穫率、すなわち、収集したウェブページの中にユーザの要求と関連する情報が含まれる割合は、最良の場合で50%程度とされており（文献1）、残りの50%のページは利用されずに捨てられることになる。文献1と文献2では、収集の効率を改善するための方法が開示されているが、そもそもインターネット上には、有用でない情報も多数含まれている。例えば、ユーザの収集要求をキーワード集合で記述した場合、そのキーワード集合を多く含んだ文書でさえ、ユーザにとって実際に有用であるとは限らず、古い情報や誤った情報、冗長な情報である可能性がある。したがって、収集効率の改善には限界があり、収集された情報が有用かどうかの判断はユーザに委ねざるを得ない。また、個々のユーザが個別にクローラを利用することは、通信ネットワークやプロキシサーバ、ウェブサーバなどにかかる負荷が大きくなるため、現実的でない。従って、より効率的な収集方法と、収集結果を無駄にせず再利用する方法が望まれる。

【0007】

さらに、クロールによってウェブページを収集するには、収集の条件として、収集を開始する起点のURLや収集する範囲、キーワードなどの条件をユーザが指定する必要が

10

20

30

40

50

ある。しかし、どのような条件を指定すれば有用な情報が得られるかが不可知である上、上述のように収集効率が良くない。従って、一般的に、検索サービスサイトや、配信型の情報フィルタリングシステムと比較して、クローラを利用するには熟練を要する。このため、有用な情報を効率よく収集するための知識やノウハウをユーザ間で共有することが望まれる。

【0008】

以上のような理由のため、クローラは、主に、検索サービスサイトが、任意の内容のウェブページを大量に収集してインデキシングする目的と、既知の限定されたウェブサイトを定期的に巡回して、更新された情報の有無を監視する目的に利用されるにとどまっている。従って、クローラが、未知の情報源から積極的に情報を収集したり、潜在的にユーザの関心に合致するであろう新しい情報を発見したりする目的に活用されていないのが現状である。

10

【0009】

一方、コミュニティのメンバーが電子掲示板等の従来のコミュニケーション手段を用いて情報をやり取りする方法では、メンバー各々の知識や専門性を生かした情報の共有を柔軟に行うことができる。しかしこの方法は、個々のユーザの能力と自発性に依存するところが大きい。有用な情報を探して他のメンバーに知らせる作業は労力を要するし、そもそも、コミュニティのメンバー全員が知らないような新しい情報を発見することは不可能である。文献3と文献4には、コミュニティでやり取りされるメッセージを解析して、ユーザの関心や嗜好（プロファイル）を求める発明が開示されているが、これらの発明は、コミュニティのメンバーの関心・嗜好に合った情報を新たに収集する手段を提供するものではない。

20

【0010】

また、有用な情報が個々のメンバーの努力によって数多く得られたとしても、その各々が未整理のまま別々のメッセージに分散している状態では、収集した情報を有効活用することができない。有用な情報を大量のメッセージの中から選び出してコミュニティのメンバー間で共有できる形に整理する作業には労力を要するが、その作業もメンバー各々の自発的な手作業に負っている。文献3に係る発明はユーザの嗜好調査、文献4に係る発明は、ユーザを関心・嗜好に基づいてカテゴライズした結果を明示することにより、コミュニケーションの円滑化を図ることを目的とする。いずれの発明も、コミュニティのメンバーのために有用な情報を整理したり保守するという作業を支援するものではない。

30

【0011】

本発明は、上記の課題を解決するためになされたものであり、ユーザの要求を満足する情報を効率よく収集するとともに、その収集結果を複数のユーザで有効に活用し、かつ、有用な情報を継続的に整理・保守する作業を支援することを目的とする。

【0012】

【課題を解決するための手段】

前記課題を解決するために、本発明に係る情報収集システムは、ユーザの要求を満足する情報を収集して提示する情報収集システムにおいて、それぞれ複数のユーザをメンバーとする複数のコミュニティを管理するコミュニティ管理手段と、各コミュニティに属するメンバーがメッセージの送受信を行うためのメッセージ送受信手段と、前記複数のコミュニティの各々で共有されている情報をユーザが閲覧するためのコミュニティ情報提示手段と、各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集要求を編集するための収集要求編集手段と、各複数のコミュニティにおいて編集された複数の収集要求のいずれかを満足する情報を情報ネットワーク上の複数の情報源から収集する情報収集手段と、前記収集した情報に基づいて前記複数の収集要求の各々に対応する収集結果を各々生成する収集結果生成手段と、各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集結果を編集するための収集結果編集手段と、を具備し、前記コミュニティ情報提示手段は、複数のコミュニティで各々作成された複数の収集結果を、当該コミュニティまたは当該コミュニティ内で送受信されるメッセージと関連付けて、当該コミュ

40

50

ニティのメンバーおよび非メンバーのユーザに提示することを特徴とする。

【0013】

本発明に係る情報収集システムの好ましい実施態様は以下のとおりである。なお、以下の各実施態様は、単独で適用しても良いし、適宜組み合わせで適用しても良い。

【0014】

(1) コミュニティのメンバーが前記メッセージ送受信手段を用いて送受信するメッセージに基づき、当該コミュニティの収集要求及び当該コミュニティの収集結果の少なくとも一方を自動的に更新すること。

【0015】

(2) コミュニティのメンバーが前記収集結果編集手段を用いて行った収集結果の編集内容に基づき、当該収集結果に対応する収集要求を更新すること。 10

【0016】

(3) コミュニティの収集結果と、当該コミュニティの収集結果に含まれる情報を重複して含む他のコミュニティの収集結果とを関連付けて提示すること。

【0017】

(4) ユーザが入力する検索条件を満足する情報を、前記情報収集手段で収集した情報の中から検索する収集情報検索手段をさらに具備し、当該収集情報検索手段は、検索された情報と、コミュニティで作成した収集結果のうち前記検索された情報を含む収集結果とを、関連付けて提示すること。

【0018】

本発明に係る情報収集方法は、ユーザの要求を満足する情報を収集して提示する情報収集方法において、各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集要求を編集し、各コミュニティにおいて編集された複数の収集要求のいずれかを満足する情報を情報ネットワーク上の複数の情報源から収集し、前記収集した情報に基づいて前記複数の収集要求の各々に対応する収集結果を各々生成し、各コミュニティに属するメンバーが共同で、当該コミュニティにおける収集結果を編集し、複数のコミュニティが各々作成された複数の収集結果を、当該コミュニティまたは当該コミュニティ内で送受信されるメッセージと関連付けて、前記複数のコミュニティの各々で共有されている情報を、当該コミュニティのメンバーおよび非メンバーのユーザに提示することを特徴とする。 20

【0019】

本発明に係る情報収集コンピュータにユーザの要求を満足する情報を収集して提示する情報収集を実行させるプログラムは、コンピュータにユーザの要求を満足する情報を収集して提示する情報収集を実行させるプログラムにおいて、各コミュニティに属するメンバーによって共同で編集された当該コミュニティにおける収集要求を入力し、各コミュニティにおいて編集された複数の収集要求のいずれかを満足する情報を情報ネットワーク上の複数の情報源から収集し、前記収集した情報に基づいて前記複数の収集要求の各々に対応する収集結果を各々生成し、各コミュニティに属するメンバーが共同で編集された当該コミュニティにおける収集結果を入力し、複数のコミュニティで各々作成された複数の収集結果を、当該コミュニティまたは当該コミュニティ内で送受信されるメッセージと関連付けて、前記複数のコミュニティの各々で共有されている情報を、当該コミュニティのメンバ 40

【0020】

【発明の実施の形態】

以下、図面を用いて本発明の実施の形態を説明する。

【0021】

図1は、本発明の一実施形態に係る情報収集システムの構成を示す図である。図1において、コミュニティ管理部1は、複数のコミュニティを管理する。すなわち、コミュニティ管理部1は、複数のコミュニティの各コミュニティにおけるメンバーである複数のユーザの情報と、各コミュニティにおけるユーザ間で送受信されるメッセージとを記憶管理する。コミュニティ管理部1は、従来技術による電子掲示板あるいはメーリングリスト等の管 50

理手段と同様に、ユーザ情報記憶部 11 とメッセージ記憶部 12 を有する。通常、コミュニティのメンバーと非メンバーでは、アクセス権、すなわち、ユーザ情報の閲覧やメッセージの送受信などが行えるか否かの権限が異なるが、このコミュニティ管理部 1 が、そのアクセス制御を行う。また、本明細書においては、ユーザは、メンバーと非メンバーを含むものとする。また、詳細は後述するように、コミュニティ管理部 1 は、ユーザからの情報収集の要求を複数記憶する収集要求記憶部 13 と、情報収集の結果としてユーザに提示する情報を複数記憶する収集結果記憶部 14 とを有する。

【0022】

コミュニティ情報提示部 2 は、複数のコミュニティの名称やメンバーなどの基本的な情報や、個々のコミュニティ内でやり取りされるメッセージや共有文書などの情報をユーザに提示する。これにより、ユーザが様々な情報を閲覧できる。

10

【0023】

メッセージ送受信部 3 は、コミュニティのメンバーが、他のメンバーに対してメッセージを送信・受信するための手段である。メッセージ送受信部 3 で送受信されたメッセージは、コミュニティ毎に整理されて、メッセージ記憶部 12 に記憶される。

【0024】

収集要求編集部 4 は、情報収集の要求をコミュニティの複数のメンバーが共同で編集して登録するための手段であり、収集要求編集部 4 で編集された結果は、収集要求記憶部 13 にコミュニティ毎に記憶される。同様に、収集結果編集部 5 は、情報収集の結果をコミュニティの複数のメンバーが利用しやすい形に編集するための手段であり、収集結果編集部 5 で編集された結果は、収集結果記憶部 14 に、コミュニティ毎に記憶される。

20

【0025】

情報収集部 6 は、収集要求記憶部 13 に記憶された複数の収集要求を入力として、インターネットやイントラネットなどの情報ネットワークから、いずれかの収集要求を満足する情報（本実施形態の場合はウェブ文書）を収集する。情報収集部 6 で収集されたウェブ文書は、ウェブ文書記憶部 7 にインデキシングされて記憶される。

【0026】

収集結果生成部 8 は、コミュニティ毎に登録された収集要求に基づき、収集したウェブ文書から要求に合致するものを選択・加工して、コミュニティ毎に収集結果を生成する。この収集結果は収集結果記憶部 14 に記憶されるが、ユーザは、必要に応じて収集結果編集部 5 を用いて収集結果をより利用しやすい形に編集して保存することができる。

30

【0027】

以上に説明した構成は、本発明を実施するための最小の構成であるが、上記の構成に加え、さらに、収集要求生成部 9 を備えてもよい。収集要求生成部 9 は、個々のコミュニティのメンバーが送受信するメッセージに基づき、当該コミュニティの収集要求を自動的に生成あるいは追加する。これと同様に、収集結果生成部 8 に、メッセージに基づいて収集結果を生成あるいは追加する機能を持たせることも可能である。さらに、収集結果生成部 8 に、ユーザが収集結果を編集した場合に、その編集内容に基づいて、対応する収集要求を変更する機能を持たせることも可能である。

【0028】

ウェブ文書検索部 10 は、情報ネットワークから収集して前記のウェブ文書記憶部 7 に記憶したウェブ文書を、ユーザが検索して利用するための手段である。ウェブ文書検索部 10 の手段の検索機能は、従来技術によるウェブ文書の検索手段と概ね同じである。本発明の実施形態に係るウェブ文書検索部 10 は、検索結果を提示する際に、収集結果記憶部 14 に記憶されている収集結果を併せて提示する処理を行う機能を備えている。

40

【0029】

以上に説明した本発明の実施形態に係る情報収集システムの構成と、従来の情報収集システムの構成との違いを、図 2 を参照して説明する。図 2 は、従来の一般的な情報収集システムの概略ブロック図である。図 2 に示す情報収集システムは、図 1 の構成要素でもある、収集要求編集部 4、収集要求記憶部 13、情報収集部 6、ウェブ文書記憶部 7、収集結

50

果生成部 8、収集結果記憶部 14、収集結果編集部 5、および、場合によりウェブ文書検索部 10を具備している。しかしながら、従来の情報収集システムは、収集要求の作成から収集結果の作成、編集までを一人のユーザが行うように構成されている。このため、従来の情報収集システムは、複数のユーザ、すなわちコミュニティで協力しあって情報を収集する目的には利用できない。また、従来の情報収集システムは、収集された情報や新たに収集すべき情報についての議論や情報交換といった活動を行うための手段も具備せず、加えて、収集結果を複数のユーザで共有し保守するための手段も備えていない。このような構成では、ユーザの労力が大きいだけでなく、複数のユーザによる情報収集結果の共有・再利用が行えないという問題がある。

【0030】

以下、本発明の実施形態について詳細に説明する。

【0031】

図3は、ユーザ情報記憶手段に記憶されるユーザの情報とコミュニティの情報を示す図である。図3(a)は、ユーザ情報31の一例であり、図3(b)は、コミュニティ情報32の一例である。ユーザ情報31は、本情報収集システムを利用する個々の登録ユーザ(所定の権限が与えられている既知のユーザ)の情報であり、ユーザID、パスワード、氏名、メールアドレス、所属コミュニティ、ホームページURL等の項目を有する。コミュニティ情報32は、本情報収集システムが管理するコミュニティの情報であり、コミュニティID、コミュニティ名、メーリングリストアドレス、掲示板URL、および、参加メンバーのユーザID等の項目を有する。メーリングリストアドレスは、コミュニティのメン

10

20

【0032】

ユーザがコミュニティを利用して情報交換の作業を行うための手順を、図4のフローチャートを用いて説明する。まず、ユーザが登録ユーザなら(ステップ41)、ユーザ認証を行う(ステップ42)。ステップ42で、認証に成功すれば(ステップ43)、当該登録ユーザの権限でコミュニティを利用できるようになる。ユーザ認証の手続きは、従来の方法と同じく、ユーザが入力したユーザのIDとパスワードを認証する方法でよい。ユーザが未登録のユーザであり、かつ新規にユーザ登録を希望するなら(ステップ44)、ユーザ登録手続き(ステップ45)を行う。ステップ45で、登録が正しく行えたならば(ステップ46)、新規の登録ユーザとしての権限でコミュニティを利用できるようになる。ユーザの登録の手続きは、従来の方法と同様に、図3(a)に示すユーザ情報31を、ユーザID31のうち氏名、パスワード等の必須項目をユーザに入力させ、ユーザIDを新しく発行することによってなされる。以上の処理は、コミュニティ管理部1が行う。

30

【0033】

その後、コミュニティ情報提示部にて、既存のコミュニティの一覧をユーザに提示する。まず、登録ユーザに対してのみ、当該ユーザが参加しているコミュニティの一覧を提示する(ステップ47)。次に、登録ユーザと、未登録の匿名ユーザの両方に対して、非参加のコミュニティの一覧を提示する(ステップ48)。ここで、コミュニティに参加していないユーザや匿名ユーザについては、制限した情報のみを提示する。図5は、複数のコミュニティ情報の一覧を提示例した画面の例を示す。この画面は、登録ユーザ「aoki」(図3のユーザIDがu1のユーザ)がユーザ認証を行った段階において提示される「aokiのポータルページ」51を示す図である。図5において、ユーザが参加しているコミュニティの一覧(図4のステップ47)を示す部分52には、「eコマース調査グループ」53や、「プロ野球ファンの集い」などの参加コミュニティが提示される。各々のコ

40

50

コミュニティに関連付けて「新着メッセージ」54、「新着情報」54の情報が提示される。新着メッセージとは、当該コミュニティに新しく送信されたメッセージの一覧である。また、新着情報とは、後述する情報収集の処理によって新しく収集された情報である。このように、コミュニティ情報の一覧画面では、コミュニティ毎に、メンバーが注目すべき新しい情報が明示的に提示される。一方、非参加コミュニティ56とは、ユーザ「aoki」が参加していないコミュニティであり、「Linuxユーザ会」57、「ガーデニング倶楽部」などのコミュニティが提示される。非参加コミュニティに関連付けて表示される「トピック」58は、当該コミュニティが関心を持って情報を収集しているトピック(話題)を表すものであり、例えば「Linuxユーザ会」57は「Linux」「ディストリビューション」などのトピックに関心を持つコミュニティであることが、非参加ユーザにも提示される一方で、具体的なメッセージなどの情報は非参加ユーザには提示されないようになっている。以上に説明したコミュニティ情報の提示の処理は図1のコミュニティ情報提示部2にて行われる。

10

【0034】

次に、ユーザが、ある一つのコミュニティを選択してこれに加わり、メッセージ送受信などの活動をおこなうための処理の流れを説明する。ユーザが、図4のステップ49にて選択したコミュニティに入ると、まず、ユーザがコミュニティのメンバーかどうかを確認される(ステップ410)。ステップ410において、ユーザが、当該コミュニティのメンバーでなく、かつ、コミュニティへの新規参加をユーザが希望するならば(ステップ411)、コミュニティへの加入手続きを行う(ステップ412)。この際、前記の登録ユーザのみがコミュニティへの加入対象となり、匿名ユーザはコミュニティへの参加は行えない。この加入手続き(ステップ412)は、図3(b)に示したコミュニティ情報32のメンバーの項目に、新規ユーザのIDを加えることによってなされるが、コミュニティの管理者や他のメンバーによって加入の可否を決定する手続きを含めてもよい。コミュニティのメンバーは、当該コミュニティ内でのメッセージの送受信と、収集要求・収集結果の閲覧・編集などの活動を行うことができる(ステップ414)。一方、コミュニティのメンバーでないユーザや、匿名ユーザは、コミュニティの利用が制限付きで許される(ステップ415)。図4に示した例では、非メンバーは、メッセージと収集結果の閲覧のみ許可され、編集は禁止するとして説明したが、コミュニティの性質によっては、これと異なるかたちで権限を許可あるいは禁止してもよい。ユーザは、コミュニティ内での活動を適宜行った後、コミュニティから出て(ステップ416)終了したり(ステップ417)、他のコミュニティに入って活動を行ってもよい。なお、図4では省略するが、コミュニティからの脱退やユーザ情報の変更、コミュニティの新規作成など、従来のコミュニティ管理システムで備えるべき処理機能も本発明の実施形態に係る情報収集システムは具備するものとする。さらに、本実施形態では主に、従来技術での電子掲示板と類似した画面例を用いて説明するが、メーリングリストのような手段を用い、ユーザ登録やコミュニティへの参加、情報の閲覧などの処理を電子メールで行うことも可能である。

20

30

【0035】

図6から図8は、コミュニティ内でのユーザの活動として、メッセージの送信、収集要求の編集、収集結果の編集の処理の流れを示した図である。また、図9から図13は、これらの処理に対応する画面の表示例である。メッセージの送受信は、本実施形態の場合は電子掲示板用の手段を用い、図1のメッセージ送受信部3にて行う。図9に示した画面例では、コミュニティ「eコマース調査グループ」に入ったユーザが、掲示板のメニュー91を選択すると、画面上には最近に送信されたメッセージ92、94、95等が表示される。個々のメッセージの間には返信関係が付されており、たとえばメッセージ95とメッセージ96は、ともにメッセージ94の返信メッセージである。画面上でユーザが一つのメッセージを選択すると、その内容が提示されるようになっており、例えば図9では、ユーザが選択したメッセージ96「著名なオークションサイト」(1月10日にユーザyama daによって送信されたメッセージ)の内容として、テキスト97が表示された場面を示している。メッセージのうち、後述する情報収集の結果と関連するものについては、こ

40

50

れらが互いに関連付けられて提示される。例えば図9では、メッセージ92「音楽配信ビジネス」に対して、このコミュニティが共同で情報収集を行っているトピック「コンテンツ配信」93が関連付けられて表示される。

【0036】

一方、メッセージの送信は、図6で示した処理により行われる。ユーザは、まずステップ61で、送信したいメッセージが既存メッセージの返信であるか、あるいは、新規のメッセージであるかを選択する。この選択は、図9の画面例ではボタン98またはボタン99を押すことによって行う。ここで、「返信メッセージ」のボタン98を押した場合は、図9で現在表示しているメッセージ96に対する返信メッセージを作成することになる。図10は返信メッセージの作成(図6のステップ62)の画面例を表す図である。ユーザは、図10に示す返信メッセージのタイトル101と本文102とを、必要ならば返信元のメッセージを引用して編集し、返信メッセージを作成する。その後、既存メッセージの返信メッセージとして送信する(ステップ63)と、上記に説明した返信関係が付されてシステムに記憶される。メッセージの送信は、図10に示す「送信」ボタン103を押すことによって行われる。新規メッセージの場合も、既存メッセージの返信メッセージの送信と同様に、図6のステップ64、ステップ65の処理を経て送信が行われる。送信されたメッセージは図1のメッセージ記憶部12に記憶され、コミュニティのメンバーが図9で説明した形で閲覧したり、新しいメッセージを返信したりすることができるようになる。

【0037】

図7は、収集要求をコミュニティのメンバーが編集する処理の流れを示す図である。収集要求は、本明細書においては、ユーザが、どのような情報を収集したいかの要求や条件を記述したデータをいい、図1の情報収集部6の入力となる。本実施形態では、収集要求をコミュニティの複数のメンバーが共同で編集することとしているため、編集内容の整合性を保つ必要がある。このため、まず、既に収集要求が存在するかどうかを確認する(ステップ71)。ステップ71において、収集要求が存在していない場合には、新規の収集要求を作成する(ステップ76)。ステップ71において、既に収集要求が存在する場合には、収集要求が他のユーザにチェックアウトされていないことを確認する(ステップ72)。この確認後に、ユーザが収集要求を編集できるようになる。ステップ72において、収集要求が他のユーザにチェックアウトされていなければ(ステップ72のYes)、まず、編集対象の収集要求が当該ユーザにチェックアウトされる(ステップ73)。そして、ユーザによる編集作業(ステップ74)の後に、チェックイン(ステップ75)を経て、システムへの登録(ステップ77)が行われる。なお、ステップ72において、収集要求が他のユーザにチェックアウトされていれば(ステップ72のNo)、当該ユーザの収集要求は編集できないので、そのまま終了する。

【0038】

以上説明した収集要求編集処理は、図1の収集要求編集部4にて行われ、編集された結果は収集要求記憶部13に記憶される。なお、編集された収集要求は、過去の収集要求と置き換えて記憶してもよいし、過去のリビジョンを保存しておいて、編集毎に新たな収集要求を追加記憶してもよい。

【0039】

図11には、収集要求を編集する画面の例を示す。ユーザが画面上で収集要求のメニュー111を選択すると、収集要求を編集するための手段が表示される。コミュニティ内で収集を行いたいトピックは、通常複数あると考えられるので、一つのコミュニティが作成する収集要求の中で、複数のトピックを記述することができるようにしている。

【0040】

図11の例では「eコマース調査グループ」というコミュニティの収集要求の例として、「電子モール」「コンテンツ配信」「オンライン・トレード」のトピックが示されている。ユーザは、これらの既存のトピックの他に新しいトピックを追加したり(ボタン116)、不要となったトピックを削除したり(ボタン113)といった編集も可能である。なお、図7で説明したチェックアウト・チェックインの処理単位は、収集要求全体を1つ

10

20

30

40

50

の処理単位とするのではなく、トピックを1つの処理単位としてもよい。個々のトピック毎に記述するデータとしては、図11に示すように、トピックの名称112、キーワード114、収集起点URL115がある。キーワード114は、収集した情報（本実施形態の場合はウェブ文書）がその内容に含むべきキーワードの論理式を記述する項目である。また、収集起点URLは、クローリングを開始するウェブ文書のURLを記述する項目である。収集起点URLは、必ずしも設定する必要はない。なぜならば、あるトピックの収集起点URLが未指定であっても、複数のコミュニティが複数のトピックに記述した収集起点URLのいずれかからクローリングすることによって、ユーザが所望する当該トピックの情報が収集できる可能性が高いからである。また、場合によっては、デフォルトの収集起点URLとして、代表的なディレクトリサイト等を選ぶことにしてもよい。以上説明した項目を図11の画面上で編集した後、「登録」ボタン117を押すことによって、編集後の収集要求がシステムに登録される。

10

【0041】

図8は、収集結果をコミュニティのメンバーが編集する処理の流れを示す図である。収集結果は、情報要求に応じてシステムが収集した情報を、コミュニティのメンバーが利用しやすい形式に加工したデータをいい、主には図1の収集結果生成部8の出力である。収集結果は、必ずしもクローリングによって収集した情報のみからなるわけではなく、ユーザが明示的に有用と思う情報を記述してもよいし、後述するように、コミュニティのメンバー間で送受信されるメッセージに含まれる情報を追加してもよい。本実施形態では、前述の収集要求と同様に、収集結果もコミュニティの複数のメンバーが共同で編集することとしているため、編集内容の整合性を保つ必要がある。このため、まず、既に収集結果が存在するかどうかを確認する（ステップ81）。ステップ81において、収集結果が存在していない場合には、新規の収集結果を作成する（ステップ86）。ステップ81において、既に収集要求が存在する場合には、収集結果が他のユーザにチェックアウトされていないことを確認する（ステップ82）。この確認後に、ユーザが編集できるようになる。ステップ82において、収集結果が他のユーザにチェックアウトされていなければ（ステップ82のYes）、まず、編集対象の収集結果がチェックアウトされる（ステップ83）。そして、ユーザによる編集作業（ステップ84）の後に、チェックイン（ステップ85）を経て、システムへの登録（ステップ87）が行われる。なお、ステップ82において、収集結果が他のユーザにチェックアウトされていれば（ステップ82のNo）、当該ユーザの収集結果は編集できないので、そのまま終了する。

20

30

【0042】

以上説明した収集結果編集処理は、図1の収集結果編集部5にて行われ、編集された結果は収集結果記憶部14に記憶される。図12には、収集結果を表示する画面の例を示す。ユーザが画面上で収集結果のメニュー121を選択すると、収集結果を表示するための手段が表示される。収集結果は、上述の収集要求のトピック毎に整理されて表示される。図12の例では、「eコマース調査グループ」の収集結果として、「電子モール」122、「コンテンツ配信」126等のトピック毎に整理されて情報が表示されている。さらに、個々のトピック中の情報は、サイト別に整理される。サイトは、インターネットにおける情報サービスの主体であり、情報源の単位でもある。図12の例では、トピック「電子モール」122の中にサイト「電子モール」123が分類されている。テキスト124は、「電子モール」123を説明するコメント文であって、コミュニティのメンバーが当該サイトの内容を理解しやすいように、メンバーの一人または複数が共同で作成したテキストである。個々のサイトの中で特に有用な情報や、新しい情報については、図12に示したように、サイト内の詳細情報125として提示する。

40

【0043】

クローリングによる情報収集の結果としては、このような既知のサイト内の情報が収集される場合（図12の情報125参照）と、新しいサイトが収集される場合（図12の情報128の例）がある。後者の場合、新しいサイトを説明するテキストはまだユーザによって作成されていないため、当該サイトのウェブ文書のテキストがそのまま提示される（

50

図12の情報129参照)が、これをより理解しやすいコメント文に直す必要がある。また一般に、クローラによって収集された情報は全てが有用な情報とは限らず、コミュニティのメンバーが共有するに値する情報を取捨・整理する作業が必要である。収集結果編集部5は、この作業をコミュニティの複数のメンバーが行うために設けられた手段であり、図13は収集結果を編集するための画面の例である。

【0044】

ユーザが図12で示した画面上の「編集」ボタン(1210)を押すと、図13に示すような画面が表示される。収集結果は上述のように、複数のトピック(「電子モール」131等)によって整理され、さらにトピックは、サイト(「電子モール」134等)によって整理される。ユーザは、新しいトピックの追加と不要なトピックの削除を行うことができる(図13のボタン1311、133)。さらに、新しいサイトの追加と不要なサイトの削除を行うことができる(図13のボタン132、136)。個々のサイト毎に編集すべき項目としては、サイト名134、サイトのURL135、サイトを説明するためのコメント文137、および、サイト内の詳細情報138である。このうち、クローリングによる情報収集で自動的に獲得できないデータはコメント文なので、ユーザの編集作業としては、コメント文を作成することが主な作業の一つであるが、これは、当該サイトのウェブ文書から取得したテキストをもとに作成すればよい。その他の作業としては、サイトや詳細情報を取捨して不要なものを削除する作業が主となる。

10

【0045】

以上の説明では、ユーザがコミュニティ内で行う活動と、そのために提供された本発明の実施形態に係る手段を中心に説明したが、以下は、ユーザが要求する情報を情報ネットワークから収集してユーザの要求に合った収集結果を生成する処理について説明する。図14は、図1の情報収集部6が行う処理の流れを表す図である。また、図14の処理の複数のステップから、収集した情報を収集結果に加える処理である図15の処理が呼び出されるが、これは図1の収集結果生成部8が行う処理である。

20

【0046】

情報収集部6は、収集対象の候補であるURLの集合を保持し、その個々のURLについて、ウェブ文書を既に取得したかどうかに係る情報や、最後に取得した日時、当該URLのリンク元URLおよびそのリンクのアンカーテキストの情報を、図1のウェブ文書記憶部7に記憶する。このURL集合をUとする。また、全コミュニティが作成する収集要求の集合をRとする。

30

【0047】

まず、Uの初期値を空集合とする(ステップ141)。その後、Rに新しい収集要求rが作成されるたびに、個々のrのトピックの収集起点URLとして新しいURLが登録されたかどうかをチェックする(ステップ142)。新しいURLu(以下、単に、「u」とのみ表記する)が登録されれば、そのスコアを計算する(ステップ143)。ここで、uのある収集要求rに対するスコアs(u,r)は、次式で計算する。

【0048】

【数1】

$$s(u,r) = \alpha * \sum_{v \rightarrow u} s(v,r) + \beta * \sum_{a: v \rightarrow u} \text{sim}(a,r) + \gamma * \text{sim}(du,r)$$

40

【0049】

ここで、 α 、 β 、 γ は定数である。vはUに含まれるURL(以下、単に、「v」とのみ表記する)であり、かつ、vはuのリンク元であるとする。s(v,r)はvの収集要求rに対するスコアである。また、a:v→uはvからuへのリンクに付されたアンカーテキストである。sim(a,r)は、アンカーテキストaと収集要求rのキーワード集合との類似度である。duはuのウェブ文書のテキストである。sim(du,r)はdu

50

のテキストと収集要求 r のキーワード集合との類似度である。収集要求 r のキーワード集合とは、収集要求 r の全てのトピックに記述されたキーワードの論理式に出現する（否定表現以外の）すべてのキーワードである。テキスト t とキーワード集合との類似度は、キーワード k の重み w_k にテキスト t 中の k の頻度 $f(t, k)$ を乗じた値を、キーワード集合の個々の要素について合計をとった値として計算する。すなわち、

【数 2】

$$\text{sim}(t, r) = \sum_{k \in r} w_k * f(t, k) / nr$$

10

とする。 nr は収集要求 r のキーワード集合の要素数である。キーワードの重み w_k は IDF (Inverted Document Frequency: すなわち、より多くのテキストに現れるキーワードほど値が小さくなる重み) で求めるのが一般的である。また、頻度 $f(t, k)$ は、単純にテキスト t 中のキーワード k の出現回数としてもよいが、テキスト t のテキスト長によって正規化した値であってもよい。 $s(u, r)$ を計算する時点で du すなわち u のウェブ文書が未取得である場合は、 $\text{sim}(du, r)$ の値は 0 とする。上記の式から分かるように、 du が未取得であっても、 u が収集要求 r を満足する可能性の大小が、 u をリンクする v のスコアや、そのリンクのアンカーテキストに基づいて推測できる。このようにして個々の収集要求 r に対する u のスコア $s(u, r)$ が求められるが、 R 中の全ての

20

収集要求 r についての $s(u, r)$ の最大値を $s(u, R)$ とする。すなわち、

$$s(u, R) = \text{Max} \{ s(u, r) \} \quad (\text{ここで、 } r \in R)$$

である。 $s(u, R)$ の値が大きい u ほど、全ての R を考慮した上で最も優先的に収集すべき URL であるとみなすことができる。

【0050】

$s(u, r)$ と $s(u, R)$ の計算方法は、上記に説明した方法に限らない。ウェブ文書が未取得の URL に対して、取得する優先順位が十分に精度良く決定できる計算方法であれば、他の計算方法を採用してもよい。優先順位の精度がよいほど、ウェブ文書を取得するコストに対して、収集要求を満たす情報が収集できる割合が高くなる。 $s(u, r)$ と $s(u, R)$ は、図 14 におけるステップ 143 とステップ 1414 のように、新たな URL に対して常に計算される。また、既知の URL に対しても、ステップ 145 とステップ 1412 のように、 R の内容が変更される毎、 u のウェブ文書や u のリンク元のスコアが変化する毎にも計算される。図 14 のステップ 144 で、ある収集要求 r のキーワードの条件が変更された場合には、ステップ 145 にて、 $s(u, r)$ と $s(u, R)$ が計算し直される。

30

【0051】

$s(u, r)$ と $s(u, R)$ をつねに最新の値に維持した上で、ステップ 146 では、URL 集合 U の中から、ウェブ文書をまだ取得していない u を選択するか、もしくは、最後にウェブ文書を取得してから閾値以上の時間が経過した URL で、かつ、スコア $s(u, R)$ が最大であるような u を選択する。そこで、 u が存在すれば (ステップ 147)、この u が、情報ネットワークから最優先に取得すべき URL である。ステップ 147 において、 u が一つも存在しなければ、取得すべき URL がないので、処理を終了する (ステップ 148) か、もしくは、収集要求集合 R の変更の有無をチェックしつつ処理を待機することになる。ステップ 149 では、 u のウェブ文書を取得する。本実施形態が対象とするインターネットのウェブ文書については、HTTP プロトコルに従った取得を行う。取得に失敗すれば (ステップ 1410)、前のステップに戻り、他の URL に対して上述の処理を繰り返し行う。取得に成功すれば、これを図 1 のウェブ文書記憶部 7 に記憶する (ステップ 1411)。次に、 u のウェブ文書の内容に基づいて、上述の $\text{sim}(du, r)$ の項を計算して、スコア $s(u, r)$ および $s(u, R)$ を計算し直す (ステップ 1412)。その後、取得したウェブ文書のパーズング (タグの解析) を行って、当該ウェブ文

40

50

書がリンクするリンク先URLを抽出し、その各々の v について(ステップ1413)、スコア $s(v, r)$ および $s(v, R)$ を計算し、URL集合 U に v を追加する(ステップ1414)。情報収集部6は、以上に説明した処理を再帰的に行い、複数のコミュニティの全ての収集要求に対して、一括して並列に、要求を満たす可能性の高いウェブ文書を収集する。したがって、個々の収集要求毎に独立にクロールを行って収集する場合と比べて、不要なウェブ文書を取得する割合が減るとともに、一つのトピックに着目したクロールでは発見しにくいような、新たな情報を発見する機会が増えるという効果がある。

【0052】

図14のステップ145、ステップ1412、及びステップ1414でスコアを計算したURLのうち、ウェブ文書を取得済みのURLの中には、個々のコミュニティの収集結果として追加すべきものがある。あるいは逆に、収集結果の中にすでに含まれているURLのうち、収集要求の条件を満たさなくなったURLについては、これを収集結果から削除する必要がある。そこで、収集結果生成部8が行う処理を図15を参照して説明する。

【0053】

まず、対象とする u のウェブ文書が取得済みであれば(ステップ151)、収集要求集合 R の中の、スコア $s(u, r)$ が変化した収集要求について、下記の処理を繰り返し行う(ステップ152)。すなわち、収集要求 r に対応する収集結果 c に既に u が含まれていれば(ステップ153)、収集要求 r の各々のトピックにキーワードの論理式の形式で記述された条件を u が満たすかどうかを調べる(ステップ154)。この処理は、 u のウェブ文書のテキストが、収集要求 r の論理式を満足する形でキーワードを含むかどうかを調べることによってなされる。 u のウェブ文書のテキストが、収集要求 r の中のどのトピックの条件も満たさなければ、 u を収集結果 c から削除する必要がある。しかし、過去にユーザが u を有用であるとみなし、収集結果 c の中に u を含めるように明示的に編集を行ったことがある場合には(ステップ155)、 u は収集結果 c から削除しない。ステップ155において、明示的な編集とは、前述の図13で示したような編集手段を用いて、 u を追加したり、あるいはコメント文などの付加情報の作成を行う編集をいう。ステップ155において、ユーザが明示的な編集を行っていない場合は、 u を収集結果 c から削除する(ステップ156)。一方、ステップ153にて、 u が収集結果 c に含まれておらず、かつ、 u が収集要求 r の条件を満たさず(ステップ157)ならば、 u は収集結果 c に追加すべきである。ただし、過去にユーザが u を不要であるとみなし、収集結果 c の中に u を含めないように明示的に編集を行ったことがある場合には(ステップ158)、 u を収集結果 c に追加しない。ステップ158において、明示的な編集とは、前述の図13で示したような編集手段を用いて u を削除した場合をいう。このような場合以外は、 u を収集結果 c に追加する(ステップ159)。ここで、本実施形態の収集結果は、図12と図13で説明したように、トピックとサイトによって整理した形式で作成されるので、 u を収集結果 c の中のトピックのうち、条件を最もよく満たすトピックの中に追加する。また、 u が既知のサイト内のURLである場合には、そのサイトの詳細情報として、図12の情報125に示したような形で追加するし、未知のサイトの情報である場合には、図12の情報128に示したように新しいサイトとして追加し、コメント文129としてウェブ文書から取得したテキストを付加する。

【0054】

本発明の実施形態に係る情報収集システムにおいては、収集要求と収集結果を、ユーザが明示的に編集するだけでなく、コミュニティ内でやり取りしたメッセージから収集要求と収集結果を自動的に更新する処理をも行う。この処理によって、動的に変化するユーザの興味・関心に常に合致するように収集要求と収集結果とを維持することができる。

【0055】

図16を用いて、メッセージに基づいて収集要求と収集結果を更新する処理の流れを説明する。

【0056】

10

20

30

40

50

未処理のメッセージ m について(ステップ161)、まず、 m の返信メッセージを再帰的に集め、 m を含むこれらのメッセージの集合を M_m とする(ステップ162)。図17に示したメッセージの例では、メッセージ171に対して、メッセージ172、173等が返信メッセージである。次に、 M_m のメッセージの各々から、URLの記述、すなわち、「http://」等で始まる記述を抽出して、これを M_m 全てのメッセージについて集めたURL集合を U_m とする(ステップ163)。図17の例では、174、176、178、1712がURLである。なお、テキスト1711は、URL174と同一であるし、メッセージ171の引用部分に含まれるので、この部分は処理しない。ステップ163の処理と同時に、 U_m の各URLに対してメッセージ中に記述されているコメント文を抽出し、 U_m の各要素に対応したコメント文集合 D_m を得る(ステップ164)。ステップ164において、メッセージからURLへのコメント文を抽出する処理は、単純には、URLと同一メッセージ内の同一の段落のテキストをそのまま抽出することで実現できるが、より複雑には、メッセージの返信関係に基づき、引用されているテキストまでも含めて文脈を理解し、複数のメッセージ間にまたがってコメント文を抽出する方法もある。図17の例では、URL174に対するテキスト175、URL176に対するテキスト177、URL178に対するテキスト179、および、URL1712に対するテキスト1711が、コメント文として抽出される。また、URL1712はURL1710(すなわち174)のサイト内のURLであり、さらに、URL1710はメッセージ171を引用した部分に含まれることから、テキスト1711およびURL1712は、URL174をより詳細に説明する情報であると解釈できる。

10

20

【0057】

このようにして、URL集合 U_m とコメント文集合 D_m とをメッセージ集合 M_m から得た後は、これを当該コミュニティの収集要求 r (または収集結果 c)の、どのトピックに追加すべきかを決定する処理を行う。

【0058】

まず、ステップ165にて、収集要求 r の各トピックに記述された収集起点URL(あるいは、収集結果 c の各トピックに記述されたURL)と、前記 U_m とを比較し、最も重複の多いトピック t_m を選択することを試みる(ステップ165)。URLの重複を調べる処理では、URLが完全に一致する場合だけでなく、URLのサイトが一致する場合も考慮する。ステップ165で t_m が選択できない場合(ステップ166)には、収集要求 r の各トピックに記述されたキーワード集合(あるいは収集結果 c の各トピックに記述されたサイト名やコメント文などのテキスト)と、 D_m のテキストとを比較し、最も重複の多いトピックを t_m とする(ステップ167)。ステップ167でも t_m が選択できない場合(ステップ168)には、トピックを新たに作成してこれを t_m とする(ステップ169)。この場合、トピック名には、メッセージのタイトルを用いる。さらに、収集要求を更新する場合には、新規トピックである t_m に対するキーワードとして、 D_m から抽出した重要語を選択する(ステップ1610)。ここでの重要語は、コメント文テキストに高い頻度で含まれ、かつ、他のトピックのコメント文テキストには低い頻度でしか含まれない語とする(従来の統計的手法により求めることができる)。ステップ165から1610の処理でトピック t_m を選択もしくは作成した後、 t_m に、先の U_m を(収集結果の更新の場合には、 D_m のコメント文と関連付けて)追加する(ステップ1611)。

30

40

【0059】

以上に説明した処理によって、図17のメッセージに対して、図18に示した収集要求、および、図19に示した収集結果が生成される。図18のトピック名181は図17のメッセージ171のタイトルであり、キーワード182は、図17のテキスト175、177、179、1711から抽出した重要語のORからなる論理式である。また、収集起点URL183には、URL174、176、178、1712が設定される。ユーザは、自動的に生成されたこれらの項目を、必要ならば前述の収集要求編集手段を用いて適宜修正して、メッセージで議論された話題に関連する情報を収集するための収集要求を簡単に作成することができる。一方、図19の収集結果については、トピック名191には図1

50

7のメッセージ171のタイトルが用いられ、サイト192、195、197にはそれぞれ図17のURL174、176、178が用いられる。各サイトに対するコメント文193、196、198には、それぞれ、図17のテキスト175、177、179が用いられる。また、メッセージ173の1711の部分は、サイト192の詳細情報として情報194に示した形で埋め込まれる。このようにして自動生成された収集結果は、常にユーザにとって利用しやすい内容に作られるとは限らず、例えばコメント文198のように余分なテキストが含まれる場合もある。この場合には、前述の収集結果編集手段を用いて、ユーザが見やすい形に自由に編集することが容易に行える。

【0060】

以上に説明した処理によって、一連のメッセージMmに対して、収集要求あるいは収集結果のトピックtmが関連付けられる(ステップ165、167)か、あるいは、新たに作成される(ステップ169)。このようなメッセージとトピックとの関連をユーザに提示することによって、ユーザがメッセージを理解したり、メッセージと関連する情報にアクセスしたりする作業を支援することができる。これは例えば、図9に示したように、メッセージ「音楽配信ビジネス」92に対して、関連するトピック「コンテンツ配信」93を関連付けて表示することによって行われる。

【0061】

一方、収集結果に対してユーザが行う編集に応じて、収集要求を自動的に更新することも可能である。この処理は、図16で説明した処理と同様の処理で実現される。ユーザが自由な形式で記述するメッセージと異なり、収集結果は、上述の収集結果編集手段(図13)で説明したような所定の形式で記述するため、この処理は図16の処理よりも比較的容易に実現できる。収集要求の条件とするキーワードは、収集結果に記述されるコメント文等から作成する。

【0062】

図1のウェブ文書検索部10の処理の流れを、図20を用いて説明する。ウェブ文書検索部10は、図1の情報収集部6が収集してウェブ文書記憶部7に記憶したウェブ文書を、ユーザが検索して利用するための手段である。

【0063】

図20において、まず、ユーザによって検索条件qが入力されると(ステップ201)、収集済みのウェブ文書からqを満足する文書を検索し、その結果のURL集合をUqとする(ステップ202)。次に、Uqの各々の要素uについて(ステップ203)、uを含む収集結果cを探す(ステップ204)。この収集結果cは、u自体を含む収集結果であってもよいし、あるいは、uと同一サイトのURLや、uをリンクするリンク元のURLを含む収集結果であってもよい。このような収集結果cが存在すれば(ステップ205)、uを説明する見出しおよび説明文として収集結果cに記述されているサイト名、コメント文のテキストを用い、uと収集結果cとを関連付けてユーザに提示する(ステップ206)。収集結果cが存在しなければ、uを説明する見出しおよび説明文として、uのウェブ文書に記述されているタイトルや本文等のテキストを用いてuをユーザに提示する(ステップ207)。

【0064】

図21は、図20で説明した処理によってユーザに提示された検索結果の画面例を示す図である。ユーザが入力した検索条件「オークション」211に対して検索された個々のウェブ文書のURL「http://xyz.com/」212等に対して、見出し「オークション」213、説明文214等を、ステップ204で求めた収集結果、例えば図19に示すサイト名192、コメント文193を用いてユーザに提示する。さらに、図21に示すように、収集結果のトピック215を収集結果と関連付けて提示する。検索結果のURLと関連する収集結果がなければ、例えば、検索結果の説明文としてウェブ文書のテキストの一部217(一般的には、冒頭部分のテキストや、検索語が出現する近傍のテキスト)を提示する。このように、ウェブ文書からそのまま得たテキストは、意味が理解し難しかったり、必ずしもそのサイトの内容を適切に表した記述でない場合がある。これ

10

20

30

40

50

に対し、説明文 2 1 4 のように、コミュニティのメンバーが収集結果の中で記述したテキストは、簡潔で理解しやすい記述である場合が多い。また、検索結果の情報に対して図 2 1 に示すように収集結果のトピックを関連付けて表示することにより、その情報がどのような分野・文脈の情報であるかが容易に理解できるようになる。さらに、ユーザは、当該トピックに含まれる他の有用な情報を利用することができる。あるトピックに関する情報を収集しているコミュニティは、そのトピックに関心を持つ専門家の集団であると言えるので、検索結果中の個々の情報について、どのようなコミュニティがこれを有用とみなしているか、いないかを、即座に知ることができるという効果もある。

【 0 0 6 5 】

以上に説明した処理は、検索結果と収集結果とを関連付けて提示する処理であったが、これと同様の方法により、あるコミュニティの収集結果に対して、他のコミュニティの収集結果を関連付けて表示することも可能である。

10

【 0 0 6 6 】

図 1 2 の情報 1 2 7 の例では、「e コマース調査グループ」が「コンテンツ配信」のトピックとして収集した情報「x x エンターテインメント」に対し、別のコミュニティである「カラオケ友の会」が収集した「家庭用コンテンツ」のトピック 1 2 7 が関連付けて提示される。この処理も、図 2 0 のステップ 2 0 4 と同様に、ある URL が収集結果に含まれているかどうかを調べることで実現される。このように、検索結果や収集結果に対し、他のコミュニティが関心のあるトピックや収集した情報を関連付けて提示することは、ユーザが検索結果や収集結果を利用する際の手助けになるだけでなく、ユーザが参加していない他のコミュニティがどのようなトピックに関心を持って活動を行っているかを、知る機会を増やす働きをする。その結果、複数のコミュニティ間の交流が活発になるという効果がある。

20

【 0 0 6 7 】

本発明は、上記の発明の実施の形態に限定されるものではない。本発明の要旨を変更しない範囲で種々変形して実施できるのは勿論である。

【 0 0 6 8 】

【発明の効果】

以上説明したように、本発明によれば、共通の関心を持ったコミュニティのメンバーが共同で収集要求と収集結果を編集し、これを継続的に洗練・保守していくことができるので、メンバー一人一人の少ない労力の寄与によって、コミュニティ全員にとって有用な情報を収集・整理して共有することができる。さらに、コミュニティ内で日常的に行われるメッセージのやり取りに基づいて、収集要求と収集結果が自動的に更新されるので、収集要求と収集結果を編集するユーザの作業が軽減するとともに、コミュニティの活動に応じて動的に変化する関心に対応した情報収集を行うことができる。

30

【図面の簡単な説明】

【図 1】 本発明の一実施形態である情報収集システムの構成を示す図。

【図 2】 従来の情報収集システムの構成の一例を表す図。

【図 3】 ユーザ情報の例を表す図。

【図 4】 ユーザの登録、認証およびコミュニティへの参加の処理の流れを表す図。

40

【図 5】 コミュニティ情報の一覧提示画面の例を表す図。

【図 6】 メッセージの送信の処理の流れを表す図。

【図 7】 収集要求の編集の処理の流れを表す図。

【図 8】 収集結果の編集の処理の流れを表す図。

【図 9】 メッセージの閲覧画面の例を表す図。

【図 1 0】 メッセージの編集画面の例を表す図。

【図 1 1】 収集要求の編集画面の例を表す図。

【図 1 2】 収集結果の閲覧画面の例を表す図。

【図 1 3】 収集結果の編集画面の例を表す図。

【図 1 4】 情報収集の処理の流れを表す図。

50

- 【図15】 収集結果の生成の処理の流れを表す図。
- 【図16】 メッセージから収集要求または収集結果を生成する処理の流れを表す図。
- 【図17】 メッセージの例を表す図。
- 【図18】 メッセージから生成された収集要求の例を表す図。
- 【図19】 メッセージから生成された収集結果の例を表す図。
- 【図20】 ウェブページ検索の処理の流れを表す図。
- 【図21】 ウェブページ検索の検索結果画面の例を表す図。

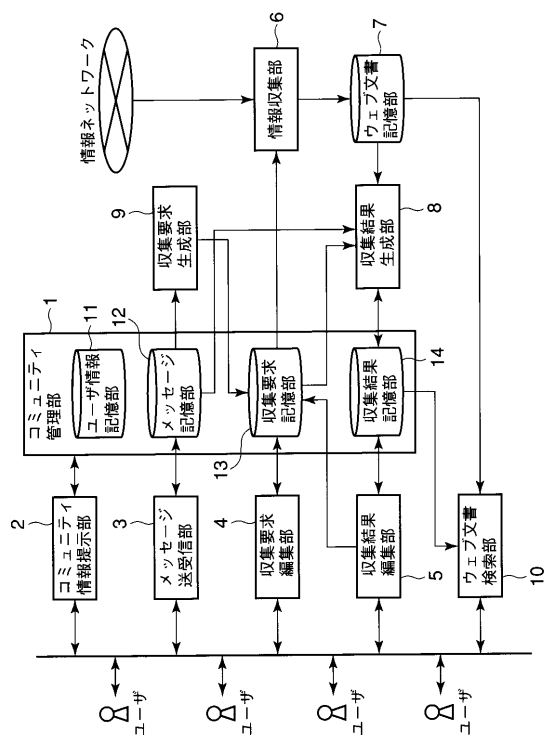
【符号の説明】

- 1 ...コミュニティ管理部
- 2 ...コミュニティ情報提示部
- 3 ...メッセージ送受信部
- 4 ...収集要求編集部
- 5 ...収集結果編集部
- 6 ...情報収集部
- 7 ...ウェブ文書記憶部
- 8 ...収集結果生成部
- 9 ...収集要求生成部
- 10 ...ウェブ文書検索部
- 11 ...ユーザ情報記憶部
- 12 ...メッセージ記憶部
- 13 ...収集要求記憶部
- 14 ...収集結果記憶部

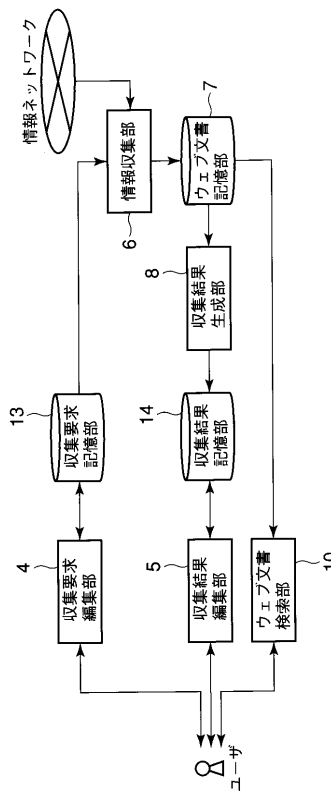
10

20

【図1】



【図2】



【 図 3 】

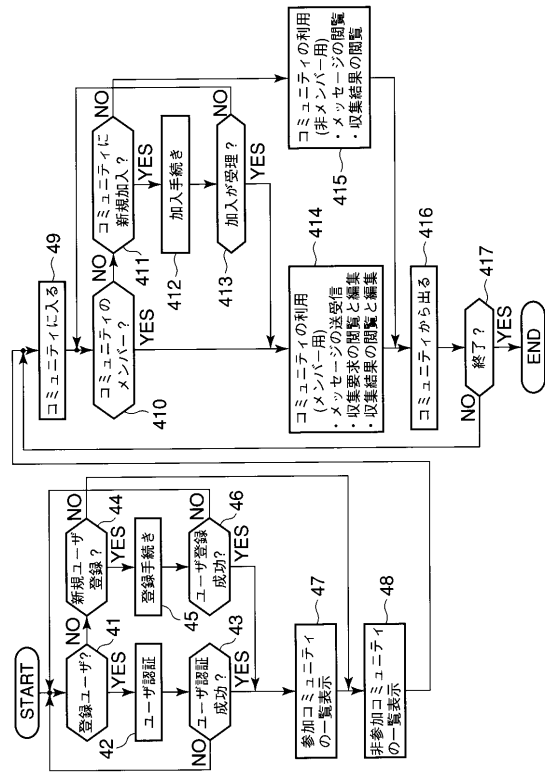
ユーザID	パスワード	氏名	メールアドレス	所属コミュニティ	ホームページURL
u1	****	青木守	aoki@aaa.net	c1,c2	http://www.aaa.net/~aoki
u2	****	伊藤隆志	itoh@bbb.com	c1,c17,c26	
...					

31 (a)

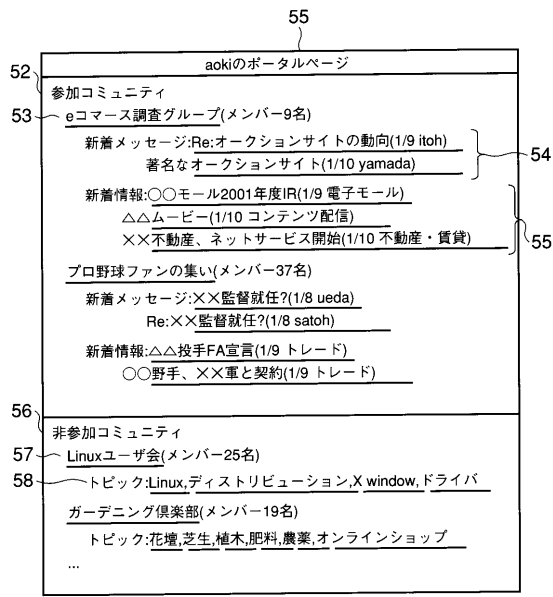
コミュニティID	コミュニティ名	メールアドレス	表示URL	メンバー
c1	eコマース調査グループ	ec-community@xxx.net	http://xxx.net/ec-community	u1,u2,u8,u36,u41...
c2	プロ野球ファンの集い	baseball-community@xxx.net		u1,u3,u19,u26,u68...
c3	Linuxユーザ会		http://xxx.net/linux-community	u5,u8,u9,u92...
...				

32 (b)

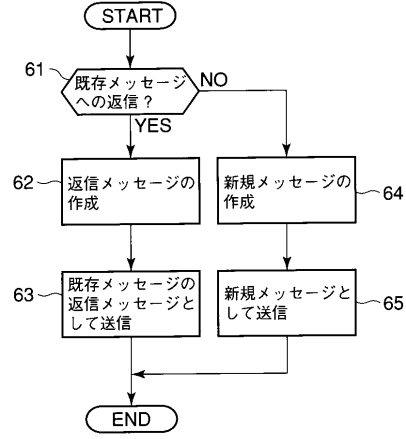
【 図 4 】



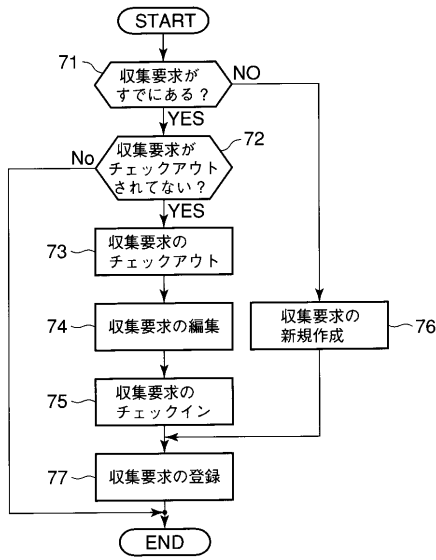
【 図 5 】



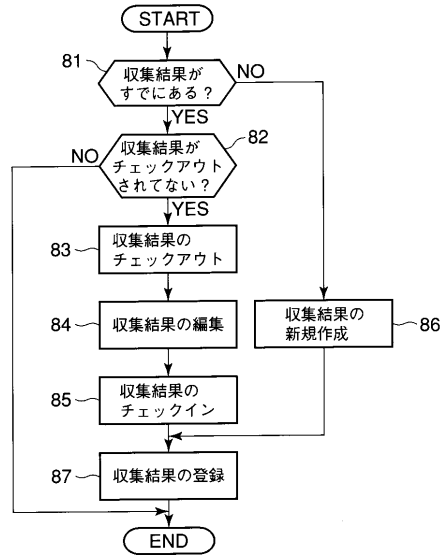
【 図 6 】



【 図 7 】



【 図 8 】



【 図 9 】

92 eコマース調査グループ

91 趣旨	掲示板
メンバー	音楽配信ビジネス(1/7 itoh) → 関連トピック
加入・脱退	Re:音楽配信ビジネス(1/8 aoki) [コンテンツ配信] 93
掲示板	オークションサイトの動向(1/9 aoki) 94
情報収集	Re:オークションサイトの動向(1/9 itoh) 95
	著名なオークションサイト(1/10 yamada) 96
	yamadaです。 著名なオークションサイトを雑誌で調べてみました。 参考とした雑誌は「〇〇マガジン12月号」と「月刊××ワールド1月号」です。
	1.〇〇オークション http://xyz.com/ オークションの草分け的サイト。ビジネスとして最初に成功した会社でもある。全世界を対象とし、出品数は1日平均5,000件。 97
	2.××オンライン http://efg.co.jp/ 日本の代表的なオークションサイト。中古車や家電製品など比較的高額な品の扱いが多い。
	...
	返信メッセージ 98
	新規メッセージ 99

【 図 10 】

eコマース調査グループ

趣旨	掲示板
メンバー	音楽配信ビジネス(1/7 itoh) → 関連トピック
加入・脱退	Re:音楽配信ビジネス(1/8 aoki) [コンテンツ配信]
掲示板	オークションサイトの動向(1/9 aoki)
情報収集	Re:オークションサイトの動向(1/9 itoh)
	著名なオークションサイト(1/10 yamada)
	返信メッセージ: Re:著名なオークションサイト 101
	aokiです。情報ありがとうございます。
	>yamadaです。 >著名なオークションサイトを雑誌で調べてみました。 >参考とした雑誌は「〇〇マガジン12月号」と「月刊××ワールド1月号」です。
	私も調べてみました。「△△マガジン11月号」からの抜粋です。 - http://hij.net/ △△ガレージセール 衣料・日用雑貨などをメインとしたオークションサイト。登録料は無料。
	...
	送信 103
	キャンセル

【図11】

図11は「eコマース調査グループ」の「情報収集要求」画面のスクリーンショットです。画面には「トピック」欄があり、「電子モール」が選択されています。キーワード欄には「(電子orオンラインor仮想) and(モールor商店街)」と入力されています。収集起点URL欄には「http://lnews.com/」と「http://e-buzz.com/」が指定されています。また、「トピックの追加」(116)、「登録」(117)、「キャンセル」(118)のボタンが配置されています。

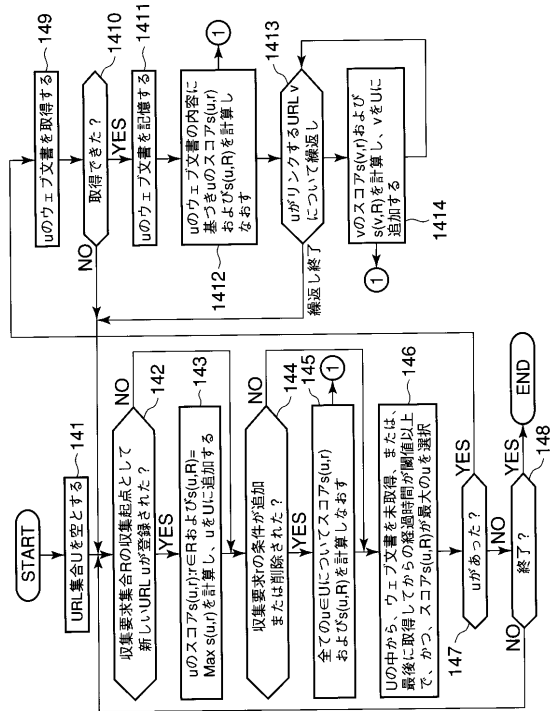
【図12】

図12は「eコマース調査グループ」の「情報収集結果」画面のスクリーンショットです。画面には「電子モール」(122)、「××市場」(123)などの収集結果がリストアップされています。各項目には「トピックの削除」(113)、「キーワード」(114)、「収集起点URL」(115)などの詳細情報が表示されています。また、「編集」(1210)のボタンも配置されています。

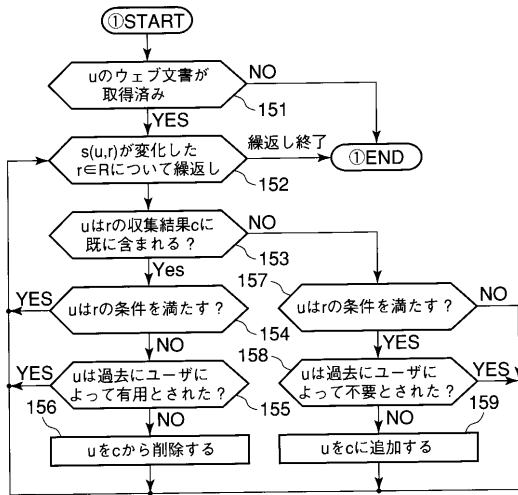
【図13】

図13は「eコマース調査グループ」の「情報収集結果(編集)」画面のスクリーンショットです。画面には「サイトの追加」(132)、「トピックの削除」(133)、「サイトの削除」(136)などの編集機能が追加されています。また、「情報の追加」(1310)のボタンも配置されています。

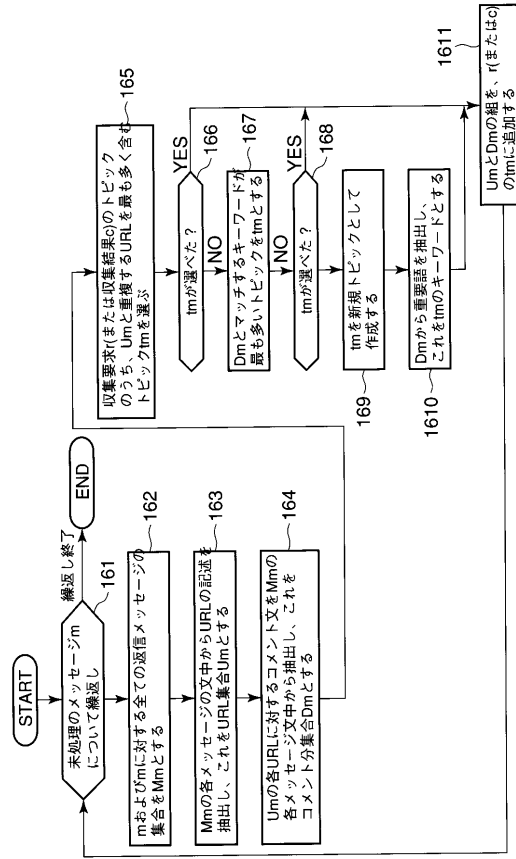
【図14】



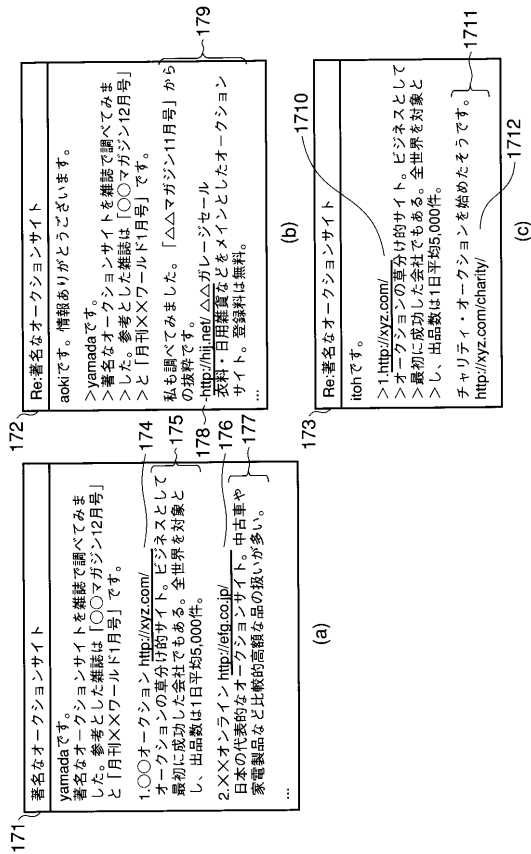
【 図 1 5 】



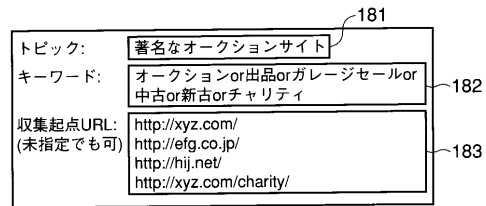
【 図 1 6 】



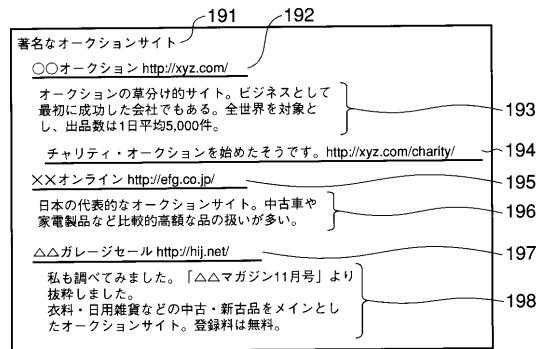
【 図 1 7 】



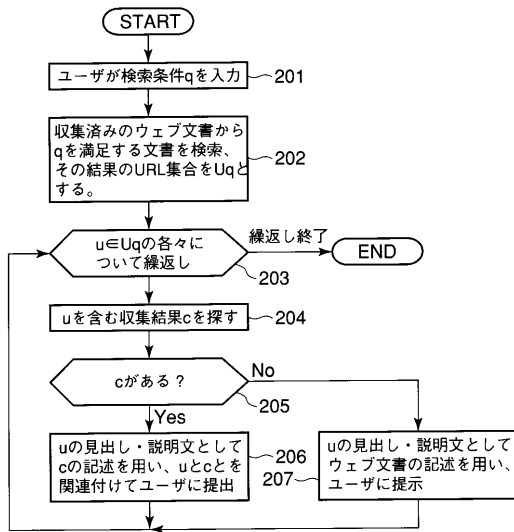
【 図 1 8 】



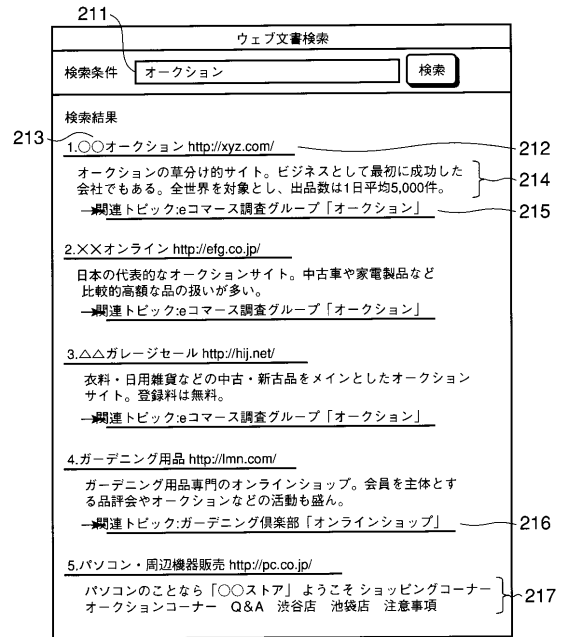
【 図 1 9 】



【 図 2 0 】



【 図 2 1 】



フロントページの続き

(72)発明者 後藤 和之

神奈川県川崎市幸区小向東芝町1番地 株式会社東芝研究開発センター内

審査官 梅本 達雄

(56)参考文献 特開2001-312515(JP,A)

特開2001-344257(JP,A)

特開2001-134616(JP,A)

イントラネット対応が進む全文検索システム,日経コンピュータ no.414 NIKKEI COMPU
TER, 1997年 3月31日

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

G06F 13/00

JSTPlus(JDream2)