

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6960274号  
(P6960274)

(45) 発行日 令和3年11月5日(2021.11.5)

(24) 登録日 令和3年10月13日(2021.10.13)

(51) Int.Cl. F 1  
G 0 6 F 1 6 / 9 5 1 ( 2 0 1 9 . 0 1 ) G O 6 F 1 6 / 9 5 1

請求項の数 9 (全 15 頁)

<p>(21) 出願番号 特願2017-160210 (P2017-160210)                  (22) 出願日 平成29年8月23日 (2017. 8. 23)                  (65) 公開番号 特開2019-40297 (P2019-40297A)                  (43) 公開日 平成31年3月14日 (2019. 3. 14)                  審査請求日 令和2年3月5日 (2020. 3. 5)</p>	<p>(73) 特許権者 319013263                  ヤフー株式会社                  東京都千代田区紀尾井町1番3号                  (74) 代理人 100149548                  弁理士 松沼 泰史                  (74) 代理人 100154852                  弁理士 酒井 太一                  (74) 代理人 100181124                  弁理士 沖田 壮男                  (74) 代理人 100194087                  弁理士 渡辺 伸一                  (72) 発明者 川崎 将平                  東京都千代田区紀尾井町1番3号 ヤフー                  株式会社内</p>
---	--

最終頁に続く

(54) 【発明の名称】 データ収集装置、データ収集方法、およびプログラム

(57) 【特許請求の範囲】

【請求項1】

ネットワークを介してアクセス可能な装置からデータを収集する収集部と、  
 前記収集部によって収集されたデータに含まれる、前記ネットワークを介してアクセス可能な装置に格納されたデータを参照するための参照情報が、所定の条件を満たす場合に、前記参照情報の一部を構成して複数の参照情報の群を特定する所属情報に対する前記収集部による収集を抑制する抑制部と

を備え、

前記所定の条件は、前記参照情報が、他の参照情報への転送を指示する情報であることである、

データ収集装置。

【請求項2】

前記所定の条件は、前記参照情報と対応するデータが、画像データまたは動画データであることである、

請求項1に記載のデータ収集装置。

【請求項3】

前記収集部により収集されたデータに基づき、前記所属情報に関するデータ収集の優先順位を決定する決定部をさらに備える、

請求項1または2に記載のデータ収集装置。

【請求項4】

前記決定部は、前記参照情報が、他の参照情報への転送を指示する情報である場合、前記他の参照情報と対応するデータに基づき、前記優先順位を決定する、  
請求項 3 に記載のデータ収集装置。

【請求項 5】

前記決定部により決定された前記所属情報に対するデータ収集の優先順位に基づき、データ収集が優先される所属情報のリストを生成する生成部をさらに備える、  
請求項 3 に記載のデータ収集装置。

【請求項 6】

前記決定部は、前記参照情報が、所定の条件を満たす場合に、前記所属情報に対するデータ収集の優先順位を下げる、  
請求項 3 に記載のデータ収集装置。

10

【請求項 7】

前記所属情報が、予め定義されたデータ収集が優先される所属情報のリストに含まれているか否かを判定し、前記所属情報が前記リストに含まれていると判定した場合、前記参照情報と対応するデータを前記収集部に収集させる判定部  
をさらに備える、  
請求項 1 から 6 のうちいずれか一項に記載のデータ収集装置。

【請求項 8】

コンピュータが、  
ネットワークを介してアクセス可能な装置からデータを収集し、  
前記収集されたデータに含まれる、前記ネットワークを介してアクセス可能な装置に格納されたデータを参照するための参照情報が、所定の条件を満たす場合に、前記参照情報の一部を構成して複数の参照情報の群を特定する所属情報に対する収集を抑制する  
データ収集方法であって、  
前記所定の条件は、前記参照情報が、他の参照情報への転送を指示する情報であること  
である、  
データ収集方法。

20

【請求項 9】

コンピュータに、  
ネットワークを介してアクセス可能な装置からデータを収集させ、  
前記収集されたデータに含まれる、前記ネットワークを介してアクセス可能な装置に格納されたデータを参照するための参照情報が、所定の条件を満たす場合に、前記参照情報の一部を構成して複数の参照情報の群を特定する所属情報に対する収集を抑制させる  
プログラムであって、  
前記所定の条件は、前記参照情報が、他の参照情報への転送を指示する情報であること  
である、  
プログラム。

30

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ収集装置、データ収集方法、およびプログラムに関する。

40

【背景技術】

【0002】

従来、ウェブから文章や画像等のデータを収集し、収集したデータを自動的にデータベース化するクローラが知られている（例えば、特許文献 1 および 2 参照）。クローラは、ウェブページ中のリンクを辿って、様々なドメインのウェブページからデータを収集する。クローラによって収集されたデータは、ウェブ情報データベースに蓄積される。ウェブ情報データベースに蓄積されたデータは、ウェブページの検索サービス等に利用される。

【先行技術文献】

【特許文献】

50

## 【 0 0 0 3 】

【特許文献 1】特開 2 0 1 2 - 6 9 1 7 1 号公報

【特許文献 2】特開平 9 - 3 2 5 9 6 8 号公報

## 【発明の概要】

【発明が解決しようとする課題】

## 【 0 0 0 4 】

検索サイト等の利便性向上のため、より多くのデータを効率的に収集することが期待されている。例えば、検索クエリに対して、単に検索クエリに対応するウェブページを提供するだけでなく、検索クエリに応じた回答を提供するシステムにおいては、より多くの有益なデータ（知識）を収集する必要がある。

10

## 【 0 0 0 5 】

しかしながら、クロールの対象となる膨大なウェブページの中から有益なデータを効率的に収集することは容易ではない。例えば、従来の幅優先探索（BFS：Breadth First Search）を用いたクロールを行う場合、不要なデータを多く含む価値の低いホストについてもクロールの対象となってしまう、リソースが無駄に消費されている場合があった。一方で、有効なデータを多く含む有益なホストに対するクロールが進まず、データ収集に時間を要してしまう場合があった。

## 【 0 0 0 6 】

本発明は、このような事情を考慮してなされたものであり、データ収集の効率を向上させることができるデータ収集装置、データ収集方法、およびプログラムを提供することを

20

【課題を解決するための手段】

## 【 0 0 0 7 】

本発明の一態様は、ネットワークを介してアクセス可能な装置からデータを収集する収集部と、前記収集部によって収集されたデータに含まれる、前記ネットワークを介してアクセス可能な装置に格納されたデータを参照するための参照情報が、所定の条件を満たす場合に、前記参照情報の一部を構成して複数の参照情報の群を特定する所属情報に対する前記収集部による収集を抑制する抑制部とを備えるデータ収集装置である。

## 【発明の効果】

## 【 0 0 0 8 】

本発明の一態様によれば、データ収集の効率を向上させることができる。

30

## 【図面の簡単な説明】

## 【 0 0 0 9 】

【図 1】実施形態のクロールサーバ 10 の使用環境を示す図である。

【図 2】実施形態のクロールサーバ 10 の構成を示す機能ブロック図である。

【図 3】実施形態のホストランク情報 D 1 の一例を示す図である。

【図 4】実施形態のコンテンツ第 1 情報 D 2 の一例を示す図である。

【図 5】実施形態のコンテンツ第 2 情報 D 3 の一例を示す図である。

【図 6】実施形態の注目ホストリスト D 4 の一例を示す図である。

【図 7】実施形態のホストランク決定部 20 の構成を示す機能ブロック図である。

40

【図 8】実施形態のクロール処理の流れの一例を示すフローチャートである。

【図 9】実施形態のホスト選択処理の流れの一例を示すフローチャートである。

【図 10】実施形態のホストランク決定処理の流れの一例を示すフローチャートである。

【図 11】実施形態の注目ホストリスト生成処理の流れの一例を示すフローチャートである。

【発明を実施するための形態】

## 【 0 0 1 0 】

## 〔概要〕

以下、図面を参照して、データ収集装置、データ収集方法、およびプログラムの実施形態について説明する。本実施形態では、データ収集装置がクロールサーバであるものとし

50

て説明する。クローラサーバとは、インターネット等のネットワークを介してアクセス可能な複数の装置からデータを自動的に収集するサーバである。本実施形態のクローラサーバは、価値の高いデータを提供する有益なホストに集中してクローラを行う。クローラサーバは、1つのプロセッサによって実現されてもよく、複数のプロセッサが分散処理することで実現されてもよい。本実施形態において、ホストとは、ネットワークを介してアクセス可能な装置に格納されたデータを参照するための参照情報（例えば、URL（Uniform Resource Locator））の一部を構成して複数の参照情報の群を特定する所属情報をいう。以下、実施形態について説明する。

#### 【0011】

##### [全体構成]

図1は、本実施形態のクローラサーバ10（データ収集装置）の使用環境を示す図である。クローラサーバ10は、画像データおよびHTML（HyperText Markup Language）データの少なくとも一方を含むページデータ（コンテンツ）を、ネットワークNWを介してアクセス可能な複数の外部サーバS1（装置）から収集する。コンテンツは、外部サーバS1に格納されており、ブラウザによって閲覧可能なページ単位のデータである。ただし、コンテンツは、ブラウザに限らず、アプリケーションプログラムによって再生されるデータでもよい。ネットワークNWは、インターネットやWAN（Wide Area Network）、LAN（Local Area Network）等を含む。

#### 【0012】

##### [クローラサーバの構成]

以下、クローラサーバ10の構成について説明する。図2は、クローラサーバ10の構成を示す機能ブロック図である。クローラサーバ10は、例えば、データ収集部12（収集部）と、解析部14（判定部）と、バッチ処理部16と、記憶部18とを備える。バッチ処理部16は、例えば、ホストランク決定部20（抑制部、決定部）と、注目ホストリスト生成部22（生成部）と、ホスト選択部24とを備える。記憶部18には、例えば、ホストランク情報D1、コンテンツ第1情報D2、コンテンツ第2情報D3、および注目ホストリストD4が記憶されている。

#### 【0013】

ホストランク情報D1には、データ収集の対象となるホストと、データ収集の優先順位を示す指標であるホストランクとが関連付けされたデータが含まれる。図3は、本実施形態のホストランク情報D1の一例を示す図である。このホストランク情報D1には、例えば、ホスト“AAA”であり、ホストランク“30”であるデータが含まれている。

#### 【0014】

コンテンツ第1情報D2には、URLと、このURLの関連情報とが関連付けされたデータが含まれる。URLの関連情報には、例えば、URLに対応するコンテンツが取得済みであるか未取得であるかを示すステータス1、コンテンツの取得の成否（例えば、HTTPステータスコード）を示すステータス2、およびデータ収集の優先度を示す指標であるスコアが含まれる。図4は、本実施形態のコンテンツ第1情報D2の一例を示す図である。このコンテンツ第1情報D2には、例えば、URLが“ddd.ddd”であり、ステータス1が“取得済”であり、ステータス2が“301”であり、スコアが“5”であるデータが含まれている。尚、初期状態のコンテンツ第1情報D2には、クローラを開始するための基礎となる複数のURLが格納されている。また、初期状態のコンテンツ第1情報D2では、全てのURLに関して、ステータス1および2は「未取得」、スコアは「未付与」となっている。

#### 【0015】

コンテンツ第2情報D3には、URLと、このURLと対応するコンテンツとが関連付けされたデータが含まれる。コンテンツには、HTMLデータおよび画像データが含まれる。図5は、本実施形態のコンテンツ第2情報D3の一例を示す図である。このコンテンツ第2情報D3には、例えば、URLが“aaa.aaa”であり、コンテンツが“HTMLデータ1”であるデータが含まれている。

10

20

30

40

50

## 【 0 0 1 6 】

注目ホストリスト D 4 には、データ収集の優先度が高いホストの一覧データが含まれる。図 6 は、本実施形態の注目ホストリスト D 4 の一例を示す図である。この注目ホストリスト D 4 には、例えば、データ収集の優先度が高いホストとして、“ B B B ”、“ E E E ”等が含まれている。

## 【 0 0 1 7 】

データ収集部 1 2 (フェッチャー) は、複数の外部サーバ S 1 からコンテンツを収集 (フェッチ) する。収集されるコンテンツには、HTML データおよび画像データが含まれる。データ収集部 1 2 は、収集したコンテンツを記憶部 1 8 のコンテンツ第 1 情報 D 2 に記憶させる。

10

## 【 0 0 1 8 】

解析部 1 4 (パーサー) は、データ収集部 1 2 により収集されて記憶部 1 8 に記憶された HTML データを解析する。例えば、解析部 1 4 は、HTML データから、ヘッダ部分を除くテキストデータを抽出し、抽出したテキストデータのなかに、新しい URL が含まれているか否かを判定する。ここで、「新しい URL」とは、記憶部 1 8 に未登録の URL である。解析部 1 4 は、抽出したテキストデータに新しい URL が含まれていると判定した場合、その URL を新しい URL として記憶部 1 8 のコンテンツ第 1 情報 D 2 (ステータス 1 および 2 は「未取得」、スコアは「未付与」) に記憶させる。

## 【 0 0 1 9 】

また、解析部 1 4 は、注目ホストリスト D 4 を参照し、上述の新しい URL に含まれるホストが注目ホストリスト D 4 に含まれているか否かを判定する。解析部 1 4 は、新しい URL に含まれるホストが注目ホストリスト D 4 に含まれていると判定した場合、このホストが優先度の高い有益なホストであると判定する。そして、解析部 1 4 は、この新しい URL を、データ収集部 1 2 の収集対象の URL のリスト (キュー) に追加する。これにより、この新しい URL に対するデータ収集が行われる。

20

## 【 0 0 2 0 】

ホストランク決定部 2 0 は、記憶部 1 8 に記憶されたコンテンツを解析し、そのコンテンツに対応する URL が属するホストに対して、データ収集の優先順位を示す指標であるホストランクを決定する。ホストランク決定部 2 0 は、日次、週次等、所定の時間間隔のバッチ処理によりホストランクを決定する。

30

## 【 0 0 2 1 】

ホストランク決定部 2 0 は、コンテンツに予め定義された特定の情報が含まれる場合に、付与する優先度を高くする。例えば、ホストランク決定部 2 0 は、コンテンツに、コンテンツの内容を示す特定の情報が含まれる場合に、付与する優先度を高くする。

## 【 0 0 2 2 】

図 7 は、本実施形態のホストランク決定部 2 0 の構成を示す機能ブロック図である。図 7 に示すように、ホストランク決定部 2 0 は、例えば、タグ情報検出部 3 0、語句検出部 3 2、スコア付与部 3 4、およびホストランク決定部 3 6 を備える。

## 【 0 0 2 3 】

タグ情報検出部 3 0 は、コンテンツのなかに、特定の情報として設定された特定のタグが含まれるか否かを検出する。「特定のタグ」は、例えば、OGP (Open Graph Protocol) タグのようなコンテンツの内容を示すテキストを含むタグである。OGP タグは、リンク先を示す URL、リンク先のコンテンツの言語、リンク先のウェブサイトの名前、リンク先のコンテンツのタイトル、リンク先のコンテンツに関する画像データの URL、リンク先のコンテンツの概要を示すテキストデータ等がひと纏まりになった情報である。

40

## 【 0 0 2 4 】

タグ情報検出部 3 0 は、コンテンツのなかに、OGP タグが含まれるか否かを検出する。例えば、タグ情報検出部 3 0 は、コンテンツのなかに OGP タグが含まれることを検出した場合、OGP タグのなかから、リンク先を示す URL、リンク先のコンテンツのタイトル、リンク先のコンテンツに関する画像データの URL、リンク先のコンテンツの概要

50

を示すテキストデータ等の情報を抽出する。また、タグ情報検出部 30 は、コンテンツのなかに OGP タグが含まれることを検出した場合、OGP タグが含まれることを示す情報と、OGP タグを含むデータに対応する URL とを対応付けてスコア付与部 34 に入力する。

【0025】

語句検出部 32 は、コンテンツのなかに、特定の情報として設定された「特定の語句」が含まれるか否かを検出する。「特定の語句」は、ウェブページのメタタグに含まれる語句であって、コンテンツの内容を示すものとして予め登録された語句でもよい。例えば、語句検出部 32 は、コンテンツに含まれるテキストデータに対して形態素解析を行い、予め登録された語句を検索することで、特定の語句が含まれるか否かを検出する。語句検出部 32 は、検出対象の特定の語句を検出した場合、特定の語句が含まれることを示す情報と、その特定の語句を含むデータに対応する URL とを対応付けてスコア付与部 34 に入力する。

10

【0026】

スコア付与部 34 は、タグ情報検出部 30 による検出結果と、語句検出部 32 による検出結果とに基づき、コンテンツに対応する URL に、データ収集の優先度を示すスコアを付与する。

【0027】

本実施形態では、スコア付与部 34 は、タグ情報検出部 30 の検出結果に基づき、URL に対して優先度として第 1 スコアを付与する。スコア付与部 34 は、タグ情報検出部 30 によってデータのなかに特定のタグが含まれることを検出した場合、データ収集の優先度が高くなるように第 1 スコアを高くする。また、本実施形態では、スコア付与部 34 は、語句検出部 32 の検出結果に基づき、コンテンツに対応する URL に対して優先度として第 2 スコアを付与する。スコア付与部 34 は、語句検出部 32 によってデータのなかに特定の語句が含まれることを検出した場合、データ収集の優先度が高くなるように第 2 スコアを高くする。

20

【0028】

ホストランク決定部 36 は、スコア付与部 34 によって付与された第 1 スコアおよび第 2 スコアの少なくとも一方に基づいて、ホストに対して、データ収集の優先順位を設定する。例えば、ホストランク決定部 36 は、URL ごとに第 1 スコアと第 2 スコアとの合計スコアを算出する。そして、ホストランク決定部 36 は、ホストごとに、このホストに属する複数の URL における合計スコアの平均値を算出する。そして、ホストランク決定部 36 は、この合計スコアの平均値が高い順に、ホストランクを決定する。

30

【0029】

また、ホストランク決定部 36 は、処理対象のホストに属する URL のなかで、所定の URL への転送（リダイレクト）を指示する URL の割合が所定の閾値以上であると判定した場合、すなわち、処理対象のホストに属する URL の多くがリダイレクトを示すものであると判定した場合、ホストランクを所定の順位だけ下げる。これにより、ホストランク決定部 36 は、コンテンツに対応する URL が所定の条件（URL の多くがリダイレクトを示すものである）を満たす場合に、この URL が属するホストに対する収集を抑制する。リダイレクト用の URL であるか否かは、データ収集部 12 によるデータ収集の際に取得した HTTP ステータスコードが、リダイレクトを示す 300 系であるか否かに基づいて判断される。

40

【0030】

また、ホストランク決定部 36 は、処理対象のホストに属する URL のなかで、コンテンツデリバリーネットワーク（CDN：Content Delivery Network）を用いて取得される画像データ、動画データ等を示す URL の割合が所定の閾値以上であると判定した場合、すなわち、処理対象のホストに属する URL の多くが画像データ等を示すものであると判定した場合、ホストランクを所定の順位だけ下げる。これにより、ホストランク決定部 36 は、コンテンツに対応する URL が所定の条件（処理対象のホストに属する URL の多く

50

が画像データ等を示すものである)を満たす場合に、このURLが属するホストに対する収集を抑制する。画像用のURLであるか否かは、URLの拡張子に基づいて判断される。画像用の拡張子には、例えば、“jpg”、“png”の拡張子が含まれる。

【0031】

また、ホストランク決定部36は、処理対象のホストに属するURLのなかで、データ収集部12によるデータ収集の際に取得したHTTPステータスコードが、サーバエラーを示す500系であると判定した場合や、認証エラーを示す400系であると判定した場合に、ホストランクを所定の順位だけ下げないようにしてもよい。また、ホストランク決定部36は、外部サーバS1からクローラを拒否する旨の情報を受け取っている場合(例えば、robot.txtに拒否URLが指定されている場合)、この拒否URLが属するホストを

10

【0032】

また、ホストランク決定部20は、処理対象のコンテンツのテキストデータに含まれるURLのリンク先のコンテンツに基づいて、ホストランクを決定してもよい。例えば、ホストランク決定部20は、処理対象のコンテンツのテキストデータに含まれるURLのリンク先のコンテンツを取得し、取得したコンテンツに対して上述の第1スコアに相当するスコア(以下、「第3スコア」という)および上述の第2スコアに相当するスコアを算出し(以下、「第4スコア」という)、第1から第4スコアに基づいて、ホストランクを決定する。

【0033】

例えば、ホストランク決定部36は、処理対象のコンテンツに対応するURLごとに、第1および第2スコアの合計スコアを算出する。さらに、ホストランク決定部36は、処理対象のコンテンツに対応するURLごとに、この処理対象のコンテンツに含まれるURLのリンク先の第3および第4スコアの合計スコアの平均値を算出する。さらに、ホストランク決定部36は、処理対象のコンテンツに対応するURLごとに、第1および第2スコアの合計スコアと、第3および第4スコアの合計スコアの平均値との2次合計スコアを算出する。そして、ホストランク決定部36は、ホストごとに、このホストに属するURLの2次合計スコアの平均値を算出し、この2次合計スコアの平均値が高い順に、ホストランキングを決定する。

20

【0034】

また、ホストランク決定部20は、処理対象のコンテンツのテキストデータに含まれるURLがリダイレクトを示すURLである場合、リダイレクト先のURLが示すコンテンツに対して、上述の第3スコアおよび第4スコアを算出し、データ収集の優先順位を設定する。

30

【0035】

注目ホストリスト生成部22は、記憶部18に記憶されたホストランク情報D1を参照し、複数のホストのなかから優先してデータを収集するホストを選出した注目ホストリストD4を生成する。例えば、注目ホストリスト生成部22は、ホストランクが所定の順位以上のホスト(例えば、上位100位)を注目ホストとして決定し、注目ホストリストD4に登録する。注目ホストリスト生成部22は、日次、週次等、所定の時間間隔のバッチ

40

【0036】

ホスト選択部24は、記憶部18に記憶されたコンテンツ第1情報D2を参照し、未だコンテンツが取得されていないURLを含むホスト(ステータス1が“未取得”であるURLが属するホスト)を選択する。さらに、ホスト選択部24は、選択したホストに属するURLのうち、コンテンツが未取得である少なくとも1つのURLをデータ収集部12の収集対象のURLのリストに追加する。これにより、キューに追加されたURLを用いたデータ収集がデータ収集部12により行われる。ホスト選択部24は、日次、週次等、所定の時間間隔のバッチ処理によりホスト選択処理を行う。

【0037】

50

クローラサーバ10の構成要素は、例えば、コンピュータにおいて、CPU (Central Processing Unit) 等のハードウェアプロセッサがプログラム (ソフトウェア) を実行することにより実現される。また、これらの構成要素のうち一部または全部は、LSI (Large Scale Integration) やASIC (Application Specific Integrated Circuit)、FPGA (Field-Programmable Gate Array)、GPU (Graphics Processing Unit) 等のハードウェア (回路部; circuitryを含む) によって実現されてもよいし、ソフトウェアとハードウェアの協働によって実現されてもよい。

【0038】

クローラサーバ10の記憶部18は、例えば、RAM (Random Access Memory)、ROM (Read Only Memory)、HDD (Hard Disk Drive)、フラッシュメモリ、またはこれらのうち複数組み合わせられたハイブリッド型記憶装置等により実現される。また、記憶部18の一部または全部は、NASや外部のストレージサーバ等、クローラサーバ10がアクセス可能な外部装置であってもよい。

10

【0039】

[クローラ処理]

以下、クローラサーバ10のクローラ処理について説明する。図8は、本実施形態のクローラ処理の流れの一例を示すフローチャートである。本フローチャートによる処理は、クローラサーバ10によって一定時間以上に亘って継続的に繰り返し実行される。尚、本フローチャートは、1つのURLを起点として実施するクローラ処理の流れを示す。

【0040】

20

まず、データ収集部12は、記憶部18に記憶されたコンテンツ第1情報D2に含まれる複数のURLのなかから、データ収集に用いるURLを選出し、選出したURLを用いてコンテンツの格納先である外部サーバS1にアクセスし、コンテンツを収集する (S101)。例えば、データ収集部12は、コンテンツ第1情報D2に含まれる複数のURLのなかから、クローラを開始するための基礎となるURLとして格納されたURL (ステータス1および2が「未取得」、スコアが「未付与」) を選出する。データ収集部12は、収集に用いたURLと、収集したコンテンツとを関連付けたデータを、記憶部18に記憶されたコンテンツ第2情報D3に追加する。

【0041】

次に、解析部14は、記憶部18に記憶されたコンテンツ第1情報D2から、データ収集部12により新たに追加されたURLとコンテンツとの組を読み出し、読み出したコンテンツを解析する (S103)。ここで、解析部14により読み出されるコンテンツは、HTMLデータである。例えば、解析部14は、読み出したHTMLデータから、ヘッダ部分を除くテキストデータを抽出し、抽出したテキストデータのなかに含まれるURLが、新しいURLであるか否かを判定する。解析部14は、抽出したテキストデータに新しいURLが含まれる場合、そのURLを新しいURL (ステータス1および2が「未取得」、スコアが「未付与」) としてコンテンツ第1情報D2に追加する (S105)。

30

【0042】

次に、解析部14は、注目ホストリストD4を参照し、上記の新しいURLが属するホストが注目ホストリストD4に含まれているか否かを判定する (S107)。解析部14は、新しいURLが属するホストが注目ホストリストD4に含まれていると判定した場合、このホストが優先度の高い有益なホストであると判定する。そして、解析部14は、この新しいURLを、データ収集部12の収集対象のURLのリスト (キュー) に追加する (S109)。これにより、データ収集部12は、キューに追加された新しいURLを用いて、再度データ収集を行い (S101)、解析部14は、再度上述の解析処理を行う (S103~S109)。これにより、優先度の高い有益なホストに属するURLが新しく発見された場合、そのURLを用いてデータを迅速に収集することができる。

40

【0043】

一方、解析部14は、新しいURLに属するホストが注目ホストリストD4に含まれていないと判定した場合、このホストが優先度の低いホストであると判定する。この場合、

50



新たなURLを用いたデータ収集は行われず、本フローチャートの処理が終了する。

【0044】

[ ホスト選択処理 ]

以下、クローラサーバ10のホスト選択処理について説明する。図9は、本実施形態のホスト選択処理の流れの一例を示すフローチャートである。本フローチャートによる処理は、日次、週次等、所定の時間間隔のバッチ処理として実行される。

【0045】

まず、ホスト選択部24は、記憶部18に記憶されたコンテンツ第1情報D2を参照し、未だコンテンツが取得されていないURLを含むホスト(ステータス1が“未取得”であるURLが属するホスト)を選択する(S201)。ここで選択されるホストは、上述のクローラ処理において、解析部14により注目ホストリストD4に含まれていないと判定され、データ収集の対象とならなかったURLが属するホストを含む。

10

【0046】

次に、ホスト選択部24は、選択したホストに属するURLのうち、コンテンツ取得済みであるURL(ステータス1が“取得済”であるURL)の数が、所定の閾値以下であるか否かを判定する(S203)。例えば、ホストごとにフェッチしたURLの数(コンテンツ取得済みURLの数)を記憶部18で管理しておき、ホスト選択部24は、このコンテンツ取得済みURLの数を参照することで、上述の判定処理を行う。ホスト選択部24は、コンテンツ取得済みであるURLの数が、閾値以下ではないと判定した場合、このホストに属するURLをデータ収集部12の収集対象に設定しない。尚、ホスト選択部24が、選択したホストに属するURLのうち、コンテンツ取得済みであるURLの割合と、所定の閾値とを比較するようにしてもよい。

20

【0047】

一方、ホスト選択部24は、コンテンツ取得済みであるURLの数が、閾値以下であると判定した場合、記憶部18に記憶されたコンテンツ第1情報D2から、選択したホストに属するURLのうち、コンテンツが未取得である(ステータス1が“未取得”である)少なくとも1つのURLを取得する(S205)。次に、ホスト選択部24は、取得したURLをデータ収集部12の収集対象のURLのリスト(キュー)に追加する(S207)。これにより、キューに追加されたURLを用いたデータ収集がデータ収集部12により行われる。

30

【0048】

ホスト選択部24は、取得したURLをデータ収集部12の収集対象のURLのリストに追加した後、または上述の判定処理においてコンテンツ取得済みであるURLの数が閾値以下ではないと判定した場合、コンテンツ第1情報D2に含まれる全てのホストに対する処理が完了したか否かを判定する(S209)。ホスト選択部24は、全てのホストに対する処理が完了していないと判定した場合、未処理のホストに対して上述のホスト選択処理(S201)以降の処理を繰り返す。一方、ホスト選択部24は、全てのホストに対する処理が完了したと判定した場合、本フローチャートの処理を終了する。

【0049】

尚、上記の実施形態においては、ホスト選択部24が、未だコンテンツが取得されていないURLに関して、コンテンツの取得処理を行わせる例を説明した。しかしながら、ホスト選択部24は、コンテンツが取得済みのURLに関して、再度、コンテンツの取得処理を行わせるようにしてもよい。これにより、コンテンツが更新された場合等に、コンテンツの最新のデータを取得することが可能である。

40

【0050】

[ ホストラंक決定処理 ]

以下、クローラサーバ10のホストラंक決定処理について説明する。図10は、本実施形態のホストラंक決定処理の流れの一例を示すフローチャートである。本フローチャートによる処理は、日次、週次等、所定の時間間隔のバッチ処理として実行される。尚、本フローチャートは、1つのホストに対するホストラंक決定処理の流れを示す。

50

## 【 0 0 5 1 】

まず、ホストランク決定部 2 0 は、記憶部 1 8 に記憶されたコンテンツ第 1 情報 D 2 から、処理対象とするホストを選択し、選択したホストに属する URL を取得する ( S 3 0 1 )。次に、ホストランク決定部 2 0 は、取得した URL に対して上述したスコア付与を行う ( S 3 0 3 )。次に、ホストランク決定部 2 0 は、その URL が属するホストのホストランクを決定する ( S 3 0 5 )。

## 【 0 0 5 2 】

次に、ホストランク決定部 2 0 は、記憶部 1 8 に記憶されたコンテンツ第 1 情報 D 2 から、取得した URL の HTTP ステータスコード ( 「ステータス 2 」 ) を取得する ( S 3 0 7 )。次に、ホストランク決定部 2 0 は、取得した URL の HTTP ステータスコードのうち、リダイレクトを示す HTTP ステータスコードの割合 ( リダイレクトを示す URL の割合 ) を算出し、この割合が所定の閾値以上であるか否かを判定する ( S 3 0 9 )。リダイレクトを示す HTTP ステータスコードは、例えば、 3 0 0 系のコードである。

## 【 0 0 5 3 】

ホストランク決定部 2 0 は、リダイレクトを示す URL の割合が所定の閾値以上であると判定した場合、すなわち、処理対象のホストに属する URL の多くがリダイレクトを示すものであると判定した場合、ホストランクを所定の順位だけ下げる ( S 3 1 1 )。一方、ホストランク決定部 2 0 は、リダイレクトを示す URL の割合が所定の閾値以上ではないと判定した場合、すなわち、処理対象のホストに属する URL にリダイレクトを示すものの数が少ないと判定した場合、上述のホストランクを下げる処理を行わない。

## 【 0 0 5 4 】

次に、ホストランク決定部 2 0 は、処理対象とするホストに属する URL のうち、画像データ、動画データ等を示す URL の割合を算出し、この割合が所定の閾値以上であるか否かを判定する ( S 3 1 3 )。ホストランク決定部 2 0 は、例えば、URL の拡張子に基づいて、URL が、画像等を示すものであるか否かを判定する。尚、ホストランク決定部 2 0 は、URL に対応するコンテンツのヘッダ情報に基づいて、URL が、画像等を示すものであるか否かを判定してもよい。

## 【 0 0 5 5 】

ホストランク決定部 2 0 は、画像等を示す URL の割合が所定の閾値以上であると判定した場合、すなわち、処理対象のホストに属する URL の多くが画像等を示すものであると判定した場合、ホストランクを所定の順位だけ下げる ( S 3 1 5 )。一方、ホストランク決定部 2 0 は、画像等を示す URL の割合が所定の閾値以上ではないと判定した場合、すなわち、処理対象のホストに属する URL に画像等を示すものの数が少ないと判定した場合、上述のホストランクを下げる処理を行わない。ホストランク決定部 2 0 は、上述の処理により決定したホストランクを記憶部 1 8 のホストランク情報 D 1 に追加または更新する。以上により、本フローチャートの処理を終了する。

## 【 0 0 5 6 】

尚、上記の実施形態においては、ホストランク決定部 2 0 が、リダイレクトを示す URL の割合が所定の閾値以上であると判定した場合や、画像等を示す URL の割合が所定の閾値以上であると判定した場合に、ホストランクを所定の順位だけ下げる例を説明した。しかしながら、ホストランク決定部 2 0 は、上述の場合に、処理対象のホストを、ホストランクから除外するようにしてもよい。

## 【 0 0 5 7 】

尚、上記の実施形態においては、リダイレクトを示す URL に対する処理と、画像等を示す URL に対する処理との両方を実施する例を説明した。しかしながら、ホストランク決定部 2 0 は、リダイレクトを示す URL に対する処理と、画像等を示す URL に対する処理とのいずれか一方を行うようにしてもよい。

## 【 0 0 5 8 】

[ 注目ホストリスト生成処理 ]

以下、クローラサーバ 1 0 の注目ホストリスト生成処理について説明する。図 1 1 は、

10

20

30

40

50

本実施形態の注目ホストリスト生成処理の流れの一例を示すフローチャートである。本フローチャートによる処理は、日次、週次等、所定の時間間隔のバッチ処理として実行される。

【0059】

まず、注目ホストリスト生成部22は、記憶部18に記憶されたホストランク情報D1を取得する(S401)。次に、注目ホストリスト生成部22は、ホストランク情報D1に含まれる複数のホストのなかから、優先してデータを収集するホストを選出した注目ホストリストD4を生成する(S403)。例えば、注目ホストリスト生成部22は、ホストランクが所定の順位以上のホスト(例えば、上位100位)を注目ホストとして決定し、注目ホストリストD4を生成する。以上により、本フローチャートの処理が終了する。

10

【0060】

以上において説明した実施形態によれば、ネットワークを介してアクセス可能な装置からデータを収集する収集部と、前記収集部によって収集されたデータに含まれる、前記ネットワークを介してアクセス可能な装置に格納されたデータを参照するための参照情報が、所定の条件を満たす場合に、前記参照情報の一部を構成して複数の参照情報の群を特定する所属情報に対する前記収集部による収集を抑制する抑制部とを備えることで、データ収集の効率を向上させることができる。すなわち、不要なデータを多く含む価値の低いホストに対するクロールを抑制し、有効なデータを多く含む有益なホストに集中してクロールを行うことができる。これにより、データ収集に要する時間を短縮し、リソースを有効に活用することができる。

20

【0061】

以上、本発明を実施するための形態について実施形態を用いて説明したが、本発明はこうした実施形態に何等限定されるものではなく、本発明の要旨を逸脱しない範囲内において種々の変形及び置換を加えることができる。

【符号の説明】

【0062】

- 10 ... クロールサーバ(データ収集装置)
- 12 ... データ収集部(収集部)
- 14 ... 解析部
- 16 ... バッチ処理部
- 18 ... 記憶部
- 20 ... ホストランク決定部
- 22 ... 注目ホストリスト生成部
- 24 ... ホスト選択部

30

【図1】

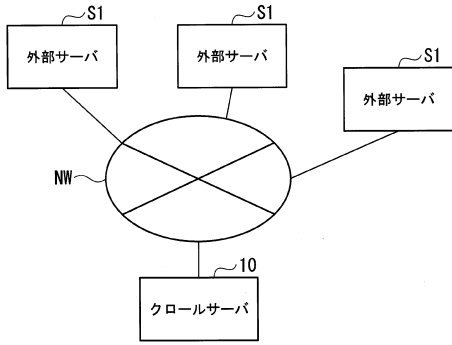


図1

【図2】

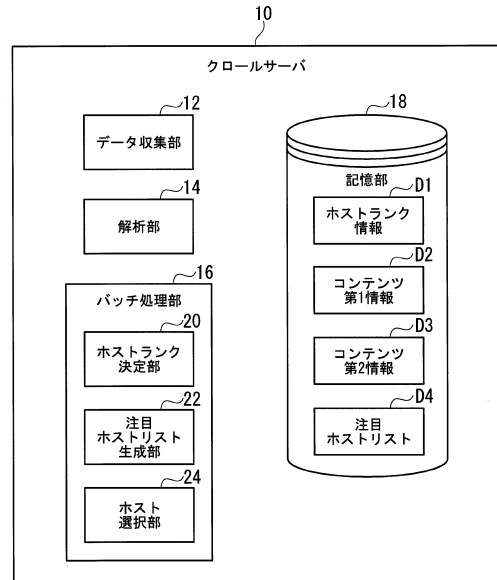


図2

【図3】

D1

ホスト	ホストランク
AAA	30
BBB	1
CCC	5
DDD	8
EEE	2

図3

【図5】

D3

URL	コンテンツ
aaa.aaa	HTMLデータ1
bbb.bbb	画像データ1
ccc.ccc	HTMLデータ5

⋮

図5

【図4】

D2

URL	ステータス1	ステータス2	スコア
ddd.ddd	取得済	301	5
eee.eee	取得済	200	70
fff.fff	未取得	未取得	未付与
ggg.ggg	取得済	401	3
hhh.hhh	取得済	200	55
iii.iii	取得済	503	40

図4

【図6】

D4

注目ホスト
BBB
EEE
QQQ
HHH
CCC
NNN

⋮

図6

【 図 7 】

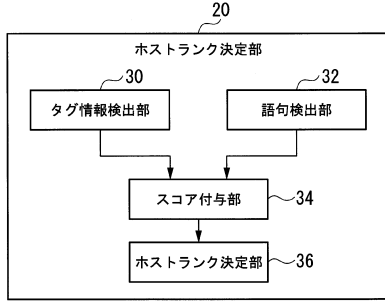


図7

【 図 8 】

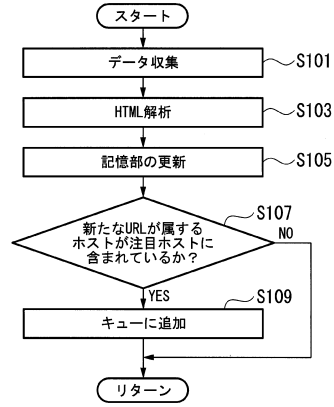


図8

【 図 9 】

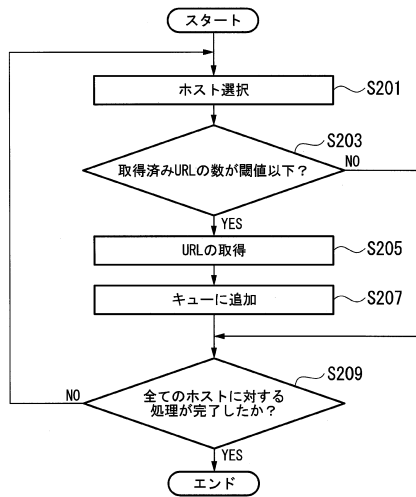


図9

【 図 10 】

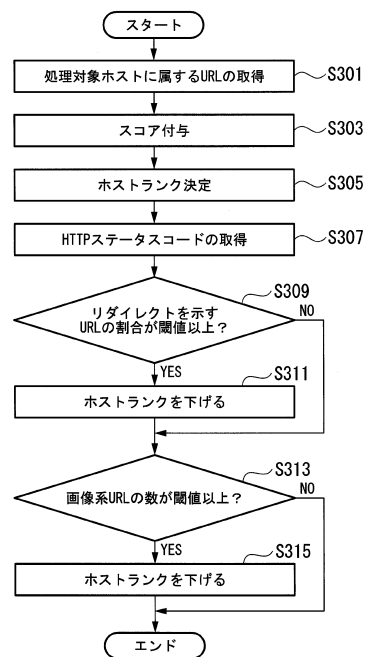


図10

【図 11】

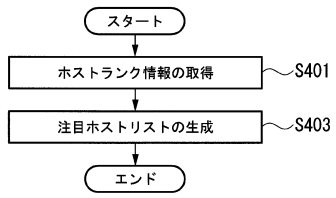


図 11

---

フロントページの続き

- (72)発明者 和良品 友大  
東京都千代田区紀尾井町1番3号 ヤフー株式会社内
- (72)発明者 俵 雄貴  
東京都千代田区紀尾井町1番3号 ヤフー株式会社内
- (72)発明者 タウフィックラチマン  
東京都千代田区紀尾井町1番3号 ヤフー株式会社内
- (72)発明者 田中 康之  
東京都千代田区紀尾井町1番3号 ヤフー株式会社内

審査官 鹿野 博嗣

- (56)参考文献 特開2006-235729(JP,A)  
特開平10-260890(JP,A)  
特開2014-186719(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G06F 16/951