

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2007-58335

(P2007-58335A)

(43) 公開日 平成19年3月8日(2007.3.8)

(51) Int. Cl.	F I	テーマコード (参考)
<b>G06F 12/00 (2006.01)</b>	G06F 12/00 546B	5B075
<b>G06F 17/30 (2006.01)</b>	G06F 17/30 419B	5B082

審査請求 未請求 請求項の数 7 O L (全 17 頁)

(21) 出願番号	特願2005-240254 (P2005-240254)	(71) 出願人	504171134 国立大学法人 筑波大学 茨城県つくば市天王台一丁目1番1
(22) 出願日	平成17年8月22日 (2005.8.22)	(74) 代理人	100091443 弁理士 西浦 ▲嗣▼晴
特許法第30条第1項適用申請有り 平成17年2月22日 電子情報通信学会データ工学専門委員会主催の「電子情報通信学会 第16回データ工学ワークショップ」のウェブサイト (http://www.digitalecity.gr.jp/~sato/DEWS2005/top04.htm) にて発表		(72) 発明者	森嶋 厚行 茨城県つくば市天王台一丁目1番1 国立大学法人筑波大学内
		(72) 発明者	飯田 敏成 茨城県つくば市天王台一丁目1番1 国立大学法人筑波大学内
		(72) 発明者	杉本 重雄 茨城県つくば市天王台一丁目1番1 国立大学法人筑波大学内

最終頁に続く

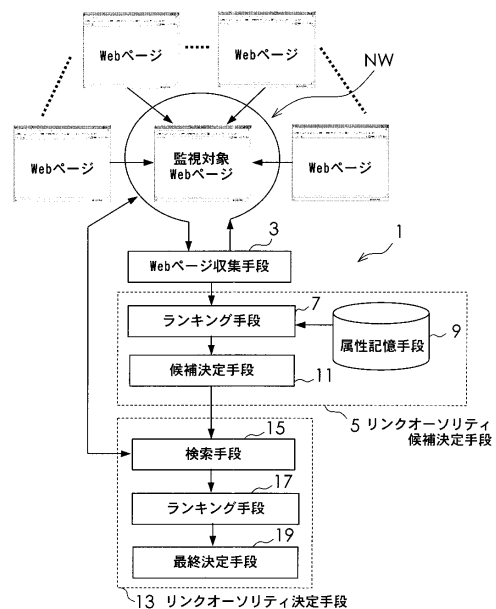
(54) 【発明の名称】 リンクオーソリティ決定方法及び装置並びにプログラム

(57) 【要約】

【課題】 できるだけ短い時間でしかも少ない労力で、適正なリンクオーソリティを決定できるリンクオーソリティの決定装置を提供する。

【解決手段】 リンクオーソリティ候補決定手段5は、Webページ収集手段3で収集した複数のWebページの中から予め定めた条件を満たす複数のWebページをリンクオーソリティ候補として定める。リンクオーソリティ決定手段13のランキング手段17は、候補となるWebページのリンクについてのリンク切れの割合が少ない順にランキングを行う。最終決定手段19は、ランキング手段17によるランキングの結果から上位のランクにある1以上のWebページをリンクオーソリティとして決定する。

【選択図】 図5



## 【特許請求の範囲】

## 【請求項 1】

監視対象とする URL へのリンクを持つ複数の Web ページを収集する収集ステップと、

前記複数の Web ページの中から予め定めた条件を満たす複数の Web ページをリンクオーソリティ候補として定めるリンクオーソリティ候補決定ステップと、

前記リンクオーソリティ候補に含まれる前記複数の Web ページの中からリンク切れを修正するために利用可能なリンクオーソリティを決定するリンクオーソリティ決定ステップとからなるリンクオーソリティの決定方法であって、

前記リンクオーソリティ候補決定ステップでは、予めリンクオーソリティとなり得る Web ページが有する複数の属性を定めて、前記複数の属性を基準にして前記複数の Web ページをリンクオーソリティとして利用可能性が高いと推測される順にランキングを行い、前記ランキングの結果から上位のランクにある複数の Web ページを前記リンクオーソリティ候補として定め、

前記リンクオーソリティ決定ステップでは、前記上位のランクにある複数の Web ページのそれぞれについて、各 Web ページについてのリンク切れのリンクの数またはリンク切れではないリンクの数を求め、この数と前記各 Web ページにあるリンクの数とに基づいてリンク切れの割合が少ない順にランキングを行い、前記ランキングの結果から上位のランクにある 1 以上の Web ページを前記リンクオーソリティとして決定することを特徴とするリンクオーソリティの決定方法。

## 【請求項 2】

前記リンクオーソリティ決定ステップでは、前記リンクの数と前記リンク切れではないリンクの数の割合を反映した値をキーとして前記ランキングを行うことを特徴とする請求項 1 に記載のリンクオーソリティの決定方法。

## 【請求項 3】

前記リンクの数と前記リンク切れではないリンクの数の割合の相乗平均を前記キーとすることを特徴とする請求項 2 に記載のリンクオーソリティの決定方法。

## 【請求項 4】

監視対象とする URL へのリンクを持つ複数の Web ページを収集する収集機能と、

前記複数の Web ページの中から予め定めた条件を満たす複数の Web ページをリンクオーソリティ候補として定めるリンクオーソリティ候補決定機能と、

前記リンクオーソリティ候補に含まれる前記複数の Web ページの中からリンク切れを修正するために利用可能なリンクオーソリティを決定するリンクオーソリティ決定機能とをコンピュータに実現させるためのプログラムであって、

前記リンクオーソリティ候補決定機能は、予め定めたリンクオーソリティとなり得る Web ページが有する複数の属性を基準にして、前記複数の Web ページをリンクオーソリティとして利用可能性が高いと推測される順にランキングを行う機能と、前記ランキングの結果から上位のランクにある複数の Web ページを前記リンクオーソリティ候補として定める機能を含み、

前記リンクオーソリティ決定機能は、前記上位のランクにある複数の Web ページのそれぞれについて、各 Web ページについてのリンク切れのリンクの数またはリンク切れではないリンクの数を求める機能と、この数と前記各 Web ページにあるリンクの数とに基づいてリンク切れの割合が少ない順にランキングを行う機能と、前記ランキングの結果から上位のランクにある 1 以上の Web ページを前記リンクオーソリティとして決定する機能とを含むことを特徴とするプログラム。

## 【請求項 5】

前記リンクオーソリティ決定機能に含まれる前記ランキングを行う機能は、前記リンクの数と前記リンク切れではないリンクの数の割合を反映した値をキーとして前記ランキングを行うことを特徴とする請求項 4 に記載のプログラム。

## 【請求項 6】

10

20

30

40

50

監視対象とするURLへのリンクを持つ複数のWebページを収集する収集手段と、前記複数のWebページの中から予め定めた条件を満たす複数のWebページをリンクオーソリティ候補として定めるリンクオーソリティ候補決定手段と、

前記リンクオーソリティ候補に含まれる前記複数のWebページの中からリンク切れを修正するために利用可能なリンクオーソリティを決定するリンクオーソリティ決定手段とからなるリンクオーソリティ決定装置であって、

前記リンクオーソリティ候補決定手段は、予め定めたリンクオーソリティとなり得るWebページが有する複数の属性を記憶する属性記憶手段と、前記属性記憶手段に記憶された前記複数の属性を基準にして、前記複数のWebページをリンクオーソリティとして利用可能性が高いと推測される順にランキングを行うランキング手段と、前記ランキング手段によるランキングの結果から上位のランクにある複数のWebページを前記リンクオーソリティ候補として定める候補決定手段を含み、

前記リンクオーソリティ決定手段は、前記上位のランクにある複数のWebページのそれぞれについて、リンク切れのリンクの数またはリンク切れではないリンクの数を求める検索手段と、前記検索手段により求めた数と前記各Webページにあるリンクの数とに基づいてリンク切れの割合が少ない順にランキングを行うランキング手段と、前記ランキング手段によるランキングの結果から上位のランクにある1以上のWebページを前記リンクオーソリティとして決定する最終決定手段とを含むことを特徴とするリンクオーソリティ決定装置。

#### 【請求項7】

前記ランキング手段は、前記リンクの数と前記リンク切れではないリンクの数の割合を反映する値をキーとして前記ランキングを行うことを特徴とする請求項6に記載のリンクオーソリティ決定装置。

#### 【発明の詳細な説明】

#### 【技術分野】

#### 【0001】

本発明は、Webページのリンク切れを自動的に修正するために用いられるリンクオーソリティを決定するために用いられるリンクオーソリティ決定方法及び装置並びにプログラムに関するものである。

#### 【背景技術】

#### 【0002】

近年、World Wide Web(以下Web)は、社会における重要なメディアの一つである。そしてWebの特徴の一つに、分散管理が挙げられる。即ち、Webコンテンツは多くの組織・個人により独立して追加・削除・更新が行われている。この特徴はWebを便利なツールとする一方で、Webコンテンツの一貫性の維持を困難としている要因でもある。コンテンツの一貫性が損なわれる一例として、Webページのリンク切れがある。そこで従来から、リンク切れが発生したときに代替りとなるリンク候補を探す技術が種々提案されている。例えば、特開平09-081446号公報「ハイパーテキストシステム」(特許文献1)には、リンク切れが起こった場合に、代替りとなるリンク先ページ候補を探すことが記載されている。この公報に記載の技術では、代替りとなるリンク先ページの発見に、アドレス(URL)の情報とWebページの内容のみを利用している。また特開平11-039327号公報「リンク情報自動修復方法および装置」(特許文献2)には、リンク切れが起こった場合に、代替りとなるリンク先ページ候補を探すために、「同じノード(ページのこと)」を探すことが記載されている。さらに、特開2001-273185号公報「ホームページアドレス登録装置及びホームページアドレス登録処理プログラムを記憶した記憶媒体」(特許文献3)には、リンク切れが起こった場合に、代替りとなるリンク先Webページ候補を探すために、Webページの内容のみを利用する技術が開示されている。

#### 【0003】

これらの技術では、リンク先ページの探索精度が必ずしも高くない。そこで本出願の発明者等は、Webのリンク切れを発見すると、変更先と考えられるリンクの候補を自動的

10

20

30

40

50

に発見しリンクの訂正を試みるシステムの開発を行い、実験を行ってきた（非特許文献1：中溝昌佳、森嶋厚行、杉本重雄及び北川博之著「WWWリンク一貫性維持支援システムにおけるリンク切れ自動修復」日本データベース学会 Letters、Vol. 3、No. 2、2004年12月。非特許文献2：中溝昌佳、森嶋厚行、有山智洋、杉本重雄及び北川博之著「WWWコンテンツ一貫性維持のためのリンク更新機構の提案」日本データベース学会 Letters、Vol. 2、No. 2、65頁 - 68頁、2003年10月）。このシステムでは、Webのリンク切れはページの移動に伴って生じたものであると仮定し移動先の探索を行う。

#### 【0004】

そして発明者等は、信頼できるリンクを含むページである「リンクオーソリティ（Link Authority）」を求めるための仕組みと、それを実装したLAサーバ（Link Authority Server）を提案し、LAサーバを利用することによりWebページの移動先を効率よく発見できる可能性があることを指摘した（非特許文献3：中溝昌佳、森嶋厚行、杉本重雄及び北川博之著「WWWにおける信頼度の高いリンクの発見」情報処理学会研究報告、Vol. 2004、No. 72（2004-DBS-134（II）、397頁 - 402頁。非特許文献4：中溝昌佳、森嶋厚行、杉本重雄及び北川博之著「WWWにおける信頼度の高いリンクの発見」電子情報通信学会技術研究報告、Vol. 104、No. 177（DE2004-63）、87頁 - 92頁、2004年7月）。

#### 【0005】

ここで発明者等が定義したWebページのリンクオーソリティとは、Webページが移動したときに、全Webページ中で十分に大きな確率でリンクが更新されるページを意味する。直観的には、リンクオーソリティとは、「リンク先の内容が変化したときに、全Webページ中で十分大きな確率でリンクが更新されるページ」のことである（GoogleなどにおけるAuthorityページとは異なる観念である）。なおここで全Webページとは、システムにおいて利用可能なWebページの全てという意味である。

#### 【0006】

例えば、図1に示す例で説明すると、ある大学Aの研究室のWebページ $u_1$ （ $m.l.s.ac.jp$ ）が存在し、このページは複数のページからリンクされているものとする。このうちWebページ $v_1$ （ $l.s.ac.jp$ ）はその研究室が所属する学科の研究室一覧ページである。このとき、一般には、Webページ $v_1$ はWebページ $u_1$ に関するリンクオーソリティである。したがって、Webページ $u_1$ がWebページ $u_2$ （ $m.org$ ）に移動したとき、Webページ $u_1$ へリンクしているページはリンク切れを起こすが、Webページ $v_1$ はWebページ $u_1$ へのリンクをWebページ $u_2$ に貼り換えるはずである。そこで発明者等は、あるWebページ $v$ は次の2つの条件を満たすとき、別のWebページ $u$ のリンクオーソリティとなるものと定義した。（1）Webページ $v$ がWebページ $u$ へのリンクを持っており、且つ（2）Webページ $u$ が別のWebページ $u_{new}$ に移動すると、Webページ $v$ 中のWebページ $u$ へのリンクがWebページ $u_{new}$ へのリンクに確実に変更されることが強く期待される。

#### 【0007】

あるWebページ中に含まれるWebページ $u$ へのリンクがリンク切れになった場合に、もしWebページ $u$ のリンクオーソリティとなるWebページ $v$ を知っていれば、Webページ $v$ を見ることにより新たなリンク先 $u_{new}$ を知ることができ、リンク先の修正が可能になる。したがってリンク先の修正をするためには、リンクオーソリティを高い精度で決定することが必要となる。

【特許文献1】特開平09-081446号公報

【特許文献2】特開平11-039327号公報

【特許文献3】特開2001-273185号公報

【非特許文献1】中溝昌佳、森嶋厚行、杉本重雄及び北川博之著「WWWリンク一貫性維持支援システムにおけるリンク切れ自動修復」日本データベース学会 Letters、V

10

20

30

40

50

o l . 3、N o . 2、2 0 0 4 年 1 2 月

【非特許文献2】中溝昌佳、森嶋厚行、有山智洋、杉本重雄及び北川博之著「WWWコンテナツ一貫性維持のためのリンク更新機構の提案」日本データベース学会 Letters、V o l . 2、N o . 2、6 5 頁 - 6 8 頁、2 0 0 3 年 1 0 月

【非特許文献3】中溝昌佳、森嶋厚行、杉本重雄及び北川博之著「WWWにおける信頼度の高いリンクの発見」情報処理学会研究報告、V o l . 2 0 0 4、N o . 7 2 ( 2 0 0 4 - D B S - 1 3 4 ( 1 1 )、3 9 7 頁 - 4 0 2 頁

【非特許文献4】中溝昌佳、森嶋厚行、杉本重雄及び北川博之著「WWWにおける信頼度の高いリンクの発見」電子情報通信学会技術研究報告、V o l . 1 0 4、N o . 1 7 7 ( D E 2 0 0 4 - 6 3 )、8 7 頁 - 9 2 頁、2 0 0 4 年 7 月

10

【発明の開示】

【発明が解決しようとする課題】

【0008】

リンクオーソリティを利用してリンク切れを修正する場合には、リンクオーソリティとしての機能ができるだけ発揮できるWebページをリンクオーソリティとして決定する必要がある。リンクオーソリティの決定方法は、種々考えられるものの、監視対象となるWebページが多くなればなるほど、決定に要するまでの時間をできるだけ短くすることができ、しかも少ない労力で実行できるものが望まれる。しかしながら従来は、この点に着目していなかったため、このような要望に答えることができる方法、装置及びプログラムはなかった。

20

【0009】

本発明の目的は、できるだけ短い時間でしかも少ない労力で、適正なリンクオーソリティを決定できるリンクオーソリティの決定方法及び装置並びにプログラムを提供することにある。

【課題を解決するための手段】

【0010】

本発明のリンクオーソリティ決定方法では、収集ステップと、リンクオーソリティ候補決定ステップと、リンクオーソリティ決定ステップとから構成される。収集ステップは、監視対象とするURLへのリンクを持つ複数のWebページを収集する。リンクオーソリティ候補決定ステップでは、収集した複数のWebページの中から予め定めた条件を満たす複数のWebページをリンクオーソリティ候補として定める。そしてリンクオーソリティ決定ステップでは、リンクオーソリティ候補に含まれる複数のWebページの中からリンク切れを修正するために利用可能なリンクオーソリティを決定する。

30

【0011】

特に本発明では、リンクオーソリティ候補決定ステップで、まず予めリンクオーソリティとなり得るWebページが有する複数の属性を定める。ここで複数の属性の定め方は任意である。しかしながらできるだけリンクオーソリティになり得るWebページが集まるように属性を定める必要がある。例えば「監視対象のWebページと同一サイトの同一ディレクトリに存在する」、「監視対象となるWebページの同一サイトの上位ディレクトリに存在する」、「監視対象となるWebページの同一サイトの上位ディレクトリに存在する」、「監視対象のWebページとの間に直接的、もしくは間接的な相互リンクが存在する」、「ファイル名にデフォルトファイル名が含まれている」等である。

40

【0012】

このような複数の属性を基準にして、複数のWebページをリンクオーソリティとして利用可能性が高いと推測される順にランキングする。ランキングの手法(何をキーとしてランキングを行うか)は任意である。そしてランキングの結果から上位のランクにある複数のWebページをリンクオーソリティ候補として定める。ここで上位のランクにある複数のWebページの定め方は任意である。例えば予め定めた順位(例えば1位及び2位)に属するWebページを候補と定める方法を採用してもよいし、また上位30件以内に入るWebページを候補と定める方法を採用することができる。

50

## 【0013】

そして本発明の方法では、リンクオーソリティ決定ステップで、各Webページについてのリンク切れのリンクの数またはリンク切れではないリンクの数を求める。そしてこの数と各Webページにあるリンクの数とに基づいてリンク切れの割合が少ない順にランキングを行い、ランキングの結果から上位のランクにある1以上のWebページをリンクオーソリティとして決定する。なお各Webページにあるリンクの数は、リンク切れのリンクの数またはリンク切れではないリンクの数を求める際に一緒に求めてもよいが、これらの数を求める場合とは、別に求めてもよい。一つのWebページにあるリンクが切れているか否かの確認作業は、時間と手間（アクセス作業、確認作業）を要する。本発明では、リンクオーソリティ候補として予め絞られた複数のWebページに関してだけ、この確認作業を行うため、収集したWebページのすべてについて確認作業を行う場合と比べて、リンクオーソリティの決定までの時間と労力を少ないものとする事ができる。特に、リンク切れの割合が少ないことを基準にしてランキングすると、より適正なリンクオーソリティ（更新率の高いリンクオーソリティ）を決定できる。この場合、リンクオーソリティ決定ステップでは、リンクの数とリンク切れではないリンクの数の割合を反映した値をキーとしてランキングを行うことが好ましい。このようにすると、ランキングが容易になる上、ランキングの精度を高いものとする事ができる。なおキーとする値は、種々の演算法を用いて演算することができる。例えば、リンクの数とリンク切れではないリンクの数の割合の相乗平均により求めた値をキーとしてランキングを行うと、ランキングの精度をより高めることができる。なおリンクオーソリティとして決定する上位のランクにある1以上のWebページの数は、例えば、リンク数等を参考にして定めればよく、リンク数が少ない場合には、上位の複数のWebページをリンクオーソリティとして用いればよい。

10

20

## 【0014】

本発明の方法をコンピュータを用いて実現する場合に用いるプログラムは、監視対象とするURLへのリンクを持つ複数のWebページを収集する収集機能と、前記複数のWebページの中から予め定めた条件を満たす複数のWebページをリンクオーソリティ候補として定めるリンクオーソリティ候補決定機能と、複数のWebページの中からリンク切れを修正するために利用可能なリンクオーソリティを決定するリンクオーソリティ決定機能とをコンピュータに実現させるためのプログラムである。特にこのプログラムでは、リンクオーソリティ候補決定機能が、予め定めたリンクオーソリティとなり得るWebページが有する複数の属性を基準にして、リンクオーソリティ候補に含まれる複数のWebページをリンクオーソリティとして利用可能性が高いと推測される順にランキングを行う機能と、ランキングの結果から上位のランクにある複数のWebページを前記リンクオーソリティ候補として定める機能を含む。そしてリンクオーソリティ決定機能が、上位のランクにある複数のWebページのそれぞれについて、リンク切れのリンクの数またはリンク切れではないリンクの数を求める機能と、この数と各Webページにあるリンクの数とに基づいてリンク切れの割合が少ない順にランキングを行う機能と、このランキングの結果から上位のランクにある1以上のWebページをリンクオーソリティとして決定する機能とを含む。

30

## 【0015】

また本発明の方法を実施するリンクオーソリティ決定装置は、監視対象とするURLへのリンクを持つ複数のWebページを収集するWebページ収集手段と、収集した複数のWebページの中から予め定めた条件を満たす複数のWebページをリンクオーソリティ候補として定めるリンクオーソリティ候補決定手段と、複数のWebページの中からリンク切れを修正するために利用可能なリンクオーソリティを決定するリンクオーソリティ決定手段とからなる。そして特に、リンクオーソリティ候補決定手段は、予め定めたリンクオーソリティとなり得るWebページが有する複数の属性を記憶する属性記憶手段と、属性記憶手段に記憶された複数の属性を基準にして、リンクオーソリティ候補に含まれる複数のWebページをリンクオーソリティとして利用可能性が高いと推測される順にランキングを行うランキング手段と、ランキング手段によるランキングの結果から上位のランク

40

50

にある複数のWebページをリンクオーソリティ候補として定める候補決定手段を含む。またリンクオーソリティ決定手段は、上位のランクにある複数のWebページのそれぞれについて、リンク切れのリンクの数またはリンク切れではないリンクの数とを求める検索手段と、検索手段により求めた数と各Webページにあるリンクの数とに基づいてリンク切れの割合が少ない順にランキングを行うランキング手段と、ランキング手段によるランキングの結果から上位のランクにある1以上のWebページをリンクオーソリティとして決定する最終決定手段とを含む。なおWebページにあるリンクの数は、検索手段により一緒に求めてもよいが、別の手段により求めてもよい。

【発明の効果】

【0016】

本発明によれば、リンクオーソリティ候補として予め絞られた複数のWebページに関してだけ、リンク切れの確認作業を行うため、収集したWebページのすべてについて確認作業を行う場合と比べて、リンクオーソリティの決定までの時間を短くして、しかも決定に要する労力を少ないものとする事ができる利点が見られる。

【発明を実施するための最良の形態】

【0017】

以下図面を参照して本発明のリンクオーソリティの決定方法及び装置の実施の形態の一例を詳細に説明する。実施の形態を説明する前に、図2を用いて本発明のリンクオーソリティ決定装置1の主要部を構成するLAサーバ(リンクオーソリティ・サーバ)を用いたリンク切れの自動修正について説明する。ここでLAサーバとは、あるWebページ(監視対象となるWebページ)のURLを「u」とした場合に、このuのリンクオーソリティと考えられる候補を決定するサーバである。実際には、リンクオーソリティを一気に求めることは困難であるため、LAサーバは複数のリンクオーソリティ候補を収集し、リンクオーソリティである可能性が高いと考えられる順にランキングしたページのURLのリスト $V = [v_1, v_2, \dots]$ を、結果として出力する。

【0018】

図3は、自動修正システムの構成(アーキテクチャ)を示す図である。リンクオーソリティ決定装置1として用いられるLAサーバ以外の構成は、既に発明者等が前述の論文で発表しているシステムである。簡単化のため、ここでは、システムが監視対象とするWebページ(リンク)はURL「u」で表されるただ一つのリンクに限定する。このシステムは監視対象としてのuがリンク切れであることを発見すると、新しいリンク先 $u_{new}$ を発見し、 $u_{new}$ に自動修正することを試みる。本システムの主要な構成要素は、対象となるリンクを監視するLIM(Link Integrity Maintenance)サーバ、移動先のページのURLである $u_{new}$ の候補集合Uを収集するチェーサー(Chaser)、Uに含まれる発見された各候補に対して「移動先らしさ」を表すスコア $score_i$ を計算するマーカー(Marker)である。これらの動作をまとめた抽象アルゴリズムを図4に示す。このアルゴリズムは、簡単に説明すると次のようになる。(1)LIMサーバはuを監視する。リンク切れを発見すると次のようにチェーサー(Chaser)とマーカー(Maker)を呼び出す。(2)チェーサー(Chaser)は移動前のWebページuのコンテンツとURLの情報wを用いて、Webサーチエンジンによる候補収集やロボットによるサイト内検索を用いた候補収集を行い、Uを作成する(4行目)。(3)Markerは各 $u_i \in U$ に対し、主に移動前のWebページと候補先のWebページとの類似度やURLの関係などに基づいてスコア $score_i$ を計算する(5~7行目)。(4)LIMサーバは、 $score_i$ を用いてUの中の $u_i$ をランキングし、リストUを計算する(8行目)。

【0019】

LAサーバを追加した場合の処理は次のようになる。チェーサー(Chaser)はリンクオーソリティを含むリンク群を、移動先リンク $u_{new}$ の候補として新たにUを追加する。マーカー(Marker)はU中の候補ページのランキングの際に、もしそのページがリンクオーソリティとされているならば高いスコアを割り当てる。具体的には、候補

10

20

30

40

50

ページが、L Aサーバが求めたランキング上位5位までのリンクオーソリティ候補からリンクされている場合、元の  $s c o r e_i$  にある定数を掛けることにより、スコアを高くする。

#### 【0020】

L I Mサーバは、監視対象の各リンクに対して、利用者が(1)リンクオーソリティを明示的に指摘する手段と、(2)自動修正を行うためのスコアの閾値を明示的に指定する手段をそれぞれ備えている。(1)では、利用者によって明示的に指定されたリンクオーソリティを参照することにより、L I Mサーバはそのリンクオーソリティに従ってリンクの自動修正を行うための閾値を指定することができる。L I Mサーバは、この閾値を見て、もし指定された閾値より候補ページのスコアが大きい場合、発見されたリンクへの自動修正を行う。また、もし指定された閾値よりも候補ページのスコアが小さい場合、移動先候補の一覧ページを生成し、その生成されたページのリンクへ自動修正を行う。

10

#### 【0021】

図5は、主として前述のL Aサーバによって構成される本発明のリンクオーソリティ決定装置1の実施の形態の構成の一例を示すブロック図である。そして図6は、リンクオーソリティ決定装置1を、コンピュータを用いて実現する場合に用いるプログラムのアルゴリズムを示すフローチャートである。

#### 【0022】

リンクオーソリティ決定装置1は、Webページ収集手段3と、リンクオーソリティ候補決定手段5と、リンクオーソリティ決定手段13とから構成される。Webページ収集手段3は、監視対象とするWebページのURL「u」へのリンクを持つ複数のWebページを、ネットワークNWを介して収集する。Webページ収集手段3は、例えば「u」へのリンクを持つページの集合を計算できるものであればどのような構成でもよい。図3の例では、チェーサー(Chaser)がこの手段の一部を構成している。収集には、例えば、クローラ(WWW巡回プログラム)を用いて収集する方法や、Webアーカイブなどを利用する方法や、Web検索エンジンを利用する方法などが考えられる。具体的なL Aサーバの実装では、次の(a)及び(b)によって収集を行うことができる。(a)Google Web Service API(商標)及びAlexs Web information Service API(商標)を用いて、uへのリンクを持つページを検索する。(b)クローラを用いて、uへのリンクを持つページがある可能性が高いと考えられる場所を探索する。すなわち(1)uと同じサイト内のページ、(2)uの属するサイトをサブドメインとして含むサイトに属するページ、(3)u中のリンクが指す先のページを探索する。したがって、必ずしもuにリンクを持つ全てのページが収集されとは限らない。

20

30

#### 【0023】

リンクオーソリティ候補決定手段5は、ランキング手段7と、属性記憶手段9と、候補決定手段11とを備えて、収集した複数のWebページの中から予め定めた条件を満たす複数のWebページをリンクオーソリティ候補として定める。属性記憶手段9は、予め定めたリンクオーソリティとなり得るWebページが有する複数の属性を記憶する。この属性については、後に詳しく説明する。そしてランキング手段7は、属性記憶手段9に記憶された複数の属性を基準にして、リンクオーソリティ候補に含まれる複数のWebページをリンクオーソリティとして利用可能性が高いと推測される順にランキングする。候補決定手段11は、ランキング手段7によるランキングの結果から、上位のランクにある複数のWebページをリンクオーソリティ候補として定める。

40

#### 【0024】

またリンクオーソリティ候補決定手段13は、検索手段15と、ランキング手段17と最終決定手段19とを備えている。この実施の形態の検索手段15は、リンクオーソリティ候補中の上位のランクにある複数のWebページのそれぞれについて、各Webページにあるリンクの数と、リンク切れのリンクの数またはリンク切れではないリンクの数とを求める。なお各Webページにあるリンクの数については、検索手段15とは別の手段で

50



求めるようにしてもよいのは勿論である。

【0025】

そしてランキング手段17は、検索手段15により求めた数（リンクの数とリンク切れのリンクの数またはリンク切れではないリンクの数）と各Webページにあるリンクの数に基づいてリンク切れの割合が少ない順にランキングを行う。さらに最終決定手段19は、ランキング手段17によるランキングの結果から、上位のランクにある1以上のWebページをリンクオーソリティとして決定する。

【0026】

ランキング手段7及びランキング手段17で行う候補のランキングが、リンクオーソリティ決定装置1の本質的な処理である。図7に示すように、リンクオーソリティ決定装置1が、 $u$ のリンクオーソリティ候補 $v_i$   $V$ をランキングし $V$ 及び $V^*$ を求める処理（1）及び（2）は、例えば次のように設計することができる。まず収集した複数のWebページ $v_i$ に対して、それぞれ図8の表に示した属性のうち、「値」の欄に「真偽」と記載した属性を有するかを判定し、ランキング手段7によりランキングを行う。これらの属性は、属性記憶手段9に記憶されている。

10

【0027】

次に図8の表に示した属性のうち、下から二つの属性「#L」と「B」に基づいてランキング手段17がランキングを行い、各 $v_i$ の「リンクオーソリティらしさ」を求める。

【0028】

図8の表に示した属性を決める際に用いるヒューリスティクス（解決法）の選択肢H1～H9について、具体的に説明する。以下の選択肢H1～H8が、直接的に図8の属性として表現されている場合もあるが、間接的に図8の属性として表現されている場合もある。

20

【0029】

H1：同一ディレクトリに $u$ と $v_i$ とが存在すれば、 $v_i$ は $u$ へのリンクを確実に更新する可能性が高いと考えられる。

【0030】

H2： $v_i$ が $u$ に対して論理的に上位の存在である場合であり、この場合には、 $v_i$ はリンクオーソリティの可能性が高い。ここで論理的に上位の存在であるとは、例えば $v_i$ が学科のページであるのに対して、 $u$ が学部のページであるといった場合である。しかし確実に更新される度合いは同一ディレクトリよりはやや劣ると考えられる。

30

【0031】

H3： $v_i$ が $u$ と同一サイトの上位ディレクトリに配置されている場合、 $v_i$ は $u$ に対して論理的に上位の存在であることが多い。

【0032】

H4： $v_i$ が $u$ に対して論理的に上位の存在である場合、Webサイトの設計の方法によっては、 $v_i$ は $u$ と同一ディレクトリ内のindex.htmlとして配置されることがある。

【0033】

H5：H1～H4より、 $v_i$ が $u$ と同一ディレクトリに配置され、且つindex.htmlである場合はリンクオーソリティである可能性が非常に高い。

40

【0034】

H6： $v_i$ と $u$ とが直接的、間接的な相互リンクを持つ場合、 $v_i$ はリンクオーソリティである可能性が高い。ここで間接相互リンクとは、 $v_i$ と $u$ とのサイトの間で、異なるページを介し互いにリンクをしているような関係を指す（図9）。

【0035】

H7：同一ディレクトリ内で $v_1$ から $v_2$ へリンクがあり、逆方向に存在しない場合、 $v_2$ が $v_1$ に対して論理的に上位の存在であるとは考えにくい。

【0036】

H8：（H7と比較して） $v_1$ と $v_2$ とが同一ディレクトリに配置され、お互いに相互

50

リンクが貼られている場合、これらの間の論理的な上位下位の関係は何ともいえない。なぜなら、「戻る」などのリンクが存在するからである。

【0037】

H9：リンク切れが多いページはリンクオーソリティとは考えにくい。このヒューリスティクスは他のヒューリスティクスH1～H8とは独立していると考えられる。

【0038】

ランキング手段7では、他とは独立していると考えられる上記H9を除いたヒューリスティックスの選択肢(H1～H8)を考慮して図10のランキング付けパターン表を作成した。このパターン表の各項目は $v_i$ の属性を現しており、各属性の説明は図8に示してある。図10において、黒丸はその属性が真であることを表し、空白は偽であることを表す。ハイフンはどちらでもよいことを表す。ランキング手段7におけるランキング処理では、収集された複数のWebページから選択されたWebページ $v_i$ が与えられると、まずこのパターン表を用いて $v_i$ がどのパターンに属するのかを判定してランキング付けする。そして候補決定手段11は、ランキング結果に基づいて上位ランク(例えばランク1～ランク5)に属するWebページをリンクオーソリティ候補として決定する。なお一般に、同じランクを持つ $v_i$ は複数存在する。

【0039】

検索手段15は、上記H9を反映させるために、候補決定手段11が決定した候補となるWebページのそれぞれについて、そのWebページに設けられたリンク数とリンク切れではないリンクの数(またはリンク切れのリンクの数)を検索する。そしてランキング手段17は、リンクの数とリンク切れではないリンクの数の割合を反映した値をキーとしてランキングを行う。

【0040】

キーとして用いる値の演算方法としては、例えば、Webページに設けられたリンク数とリンク切れではないリンクの数との割合の相乗平均

【数1】

$$\sqrt{\#L \times B}$$

【0041】

を求めることが考えられる。本実施の形態では、この相乗平均をキーとして降順に並べるランキングを実行する。そしてランキングにより並べられた結果を $V^*$ とする。最終決定手段19は、この結果 $V^*$ の中が上位に位置する1以上のWebページをリンクオーソリティとして決定する。

【0042】

なお $v_i$ が与えられたとき、図8に示した相互リンクの属性を求めるためには、 $u$ と $v_i$ との間に相互リンクがあるか否かを判断しなければならない。直接的な相互リンクを持つことは簡単に調べられるが、図9のような間接的な相互リンクを発見するためには多くのリンクを探索する必要がある。本来、この処理では図9の $site_u$ と $site_v$ との両サイトの内部を探索しなければならない。しかし予め処理(1)によって $v_i$ から $u$ へリンクの存在が保証されていることを利用する。具体的には $site_u$ だけの探索を行い、 $site_v$ 中のいずれかのページに対するリンクを発見すると、 $u$ と $v_i$ との間に間接的な相互リンクが存在するとみなす。これは厳密にいえば $site_v$ の $v'$ と $v_i$ 間のリンクの存在を保証しない近似的な処理であるが、これにより探索処理を半分に行うことができる。

【0043】

次に図6に示したフローチャートを参照して本発明のリンクオーソリティの決定方法の実施の形態を、コンピュータを用いて実行する場合について説明する。本発明のリンクオーソリティ決定方法では、収集ステップ(ST1, ST2)と、リンクオーソリティ候補

10

20

30

40

50

決定ステップ ( S T 3 , S T 4 ) と、リンクオーソリティ決定ステップ ( S T 5 ~ S T 1 0 ) とから構成される。収集ステップでは、監視対象とする複数の W e b ページから 1 つの W e b ページを選択し ( S T 1 )、この W e b ページを監視対象として、この W e b ページの U R L へのリンクを持つ複数の W e b ページを収集する ( S T 2 )。次にリンクオーソリティ候補決定ステップでは、収集した複数の W e b ページの中から予め定めた条件を満たす複数の W e b ページをリンクオーソリティ候補として定める。そのために、まず予めリンクオーソリティとなり得る W e b ページが有する複数の属性を定める。実際には事前に定めて記憶した属性を用いる。次に複数の属性を基準にして、複数の W e b ページをリンクオーソリティとして利用可能性が高いと推測される順にランキングする ( S T 3 )。ランキングの手法は任意である。そしてランキングの結果から上位のランクにある複数の W e b ページをリンクオーソリティ候補として定める ( S T 4 )。ここで上位のランクにある複数の W e b ページの定め方は任意であり、例えば予め定めた順位 ( 例えば 1 位及び 2 位 ) に属する W e b ページを候補と定める方法を採用してもよいし、また上位 3 0 件以内に入る W e b ページを候補と定める方法を採用することができる。

#### 【 0 0 4 4 】

次に、リンクオーソリティ決定ステップでは、リンクオーソリティ候補に含まれる複数の W e b ページの中からリンク切れを修正するために利用可能なリンクオーソリティを決定する。まず上位のランクにある複数の W e b ページのそれぞれについて、リンク切れのリンクの数またはリンク切れではないリンクの数を求める ( S T 5 , S T 6 )。そしてこの数と各 W e b ページにあるリンクの数とに基づいてリンク切れの割合が少ない順にランキングを行い ( S T 9 )、ランキングの結果から上位のランクにある 1 以上の W e b ページをリンクオーソリティとして決定する ( S T 1 0 )。なお本実施の形態では、各 W e b ページにあるリンクの数は、リンク切れのリンクの数またはリンク切れではないリンクの数を求めるときに一緒に求めている。

#### 【 0 0 4 5 】

ランキングのために、具体的には、各 W e b ページにあるリンクの数とリンク切れではないリンクの数の割合を反映した値 ( 本実施の形態では相乗平均の演算値 ) を演算する ( S T 7 )。すべての候補について、演算を実行し、その後この値 ( 相乗平均の演算値 ) をキーとしてランキングを行う ( S T 9 )。そしてランキングの上位にある 1 以上の W e b ページをリンクオーソリティとして決定する ( S T 1 0 )。リンクオーソリティとして決定する上位のランクにある 1 以上の W e b ページの数は、例えば、リンク数等を参考にし定めればよく、リンク数が少ない場合には、上位の複数の W e b ページをリンクオーソリティとして用いればよい。一つの W e b ページにあるリンクが切れているか否かの確認作業は、時間と手間 ( アクセス作業、確認作業 ) を要する。そこで本実施の形態では、リンクオーソリティ候補として予め絞られた複数の W e b ページに関してだけ、この確認作業を行う。そのため、収集した W e b ページのすべてについて確認作業を行う場合と比べて、本発明によれば、リンクオーソリティの決定までの時間と労力を少ないものとしてすることができる。特に、リンク切れの割合が少ないことを基準にしてランキングすると、より適正なリンクオーソリティ ( 更新率の高いリンクオーソリティ ) を決定できる。図 6 のアルゴリズムからなるプログラムにおいて、ステップ S T 2 によって W e b ページ収集機能が実現され、ステップ S T 3 によりランキング機能が実現され、ステップ S T 4 によってリンクオーソリティ候補を定める機能が実現され、ステップ S T 5 及び S T 6 によってリンクの数を求める機能が実現され、ステップ S T 7 からステップ法 S T 9 によってランキング機能が実現され、ステップ S T 1 0 でリンクオーソリティを決定する機能が実現されている。

#### 【 0 0 4 6 】

次に、本発明の実施の形態を用いて実験を行った結果について説明する。リンクオーソリティ決定装置によって発見されたリンクオーソリティを利用することにより、リンク切れ自動修正システムの移動先発見精度はどのように変わるか検証を行った。

この実験ではまず、筑波大学、芝浦工業大学、北海道大学、東北大学、東京大学、名古屋

屋大学、京都大学、大阪大学、九州大学の計9大学のドメインに属するサイトの中に含まれるリンクを収集した。本実験で監視対象とするリンクは、これらの収集したリンクのうち、リンク元とリンク先が異なるサイトであるようなリンク(合計49750個)である。したがって、リンクが指している先は学内のサイトとは限らない。監視対象をこれらのリンクに絞った理由は、リンク先とリンク元とが異なるサイトであればリンク切れが発生する確率が高いと考えたためである。

#### 【0047】

これらのリンクを対象に2005年1月9日より実験を行った。2005年2月4日時点で、LIMサーバは監視対象のうち146個のリンク切れを発見した(図11)。このうちWebページ移動によって生じたと考えられるリンク切れは47個存在した。ここで、ページ移動によるリンク切れとの判断は次のように行った。つまり、LIMサーバ探索ログなどを基に、様々な方法で移動先の探索を行い、移動先と考えられるページが発見できたものを移動ページとした。ここで移動先発見が成功した場合は、LIMサーバが出力した結果の上位3位以内に正しい移動先が含まれている場合とした。

10

#### 【0048】

まず、移動先の探索結果を図12に示す。図における「LAサーバなし」の行はLAサーバを利用せず、「LAサーバあり」の行はLAサーバを利用して探索を行った結果である。LAサーバを利用せずに移動先の発見に成功したものは28個、失敗したものは19個、成功率は59.6%であった。それに対してLAサーバを利用した場合は、移動先の発見に成功したものは34個、失敗したものは13個、成功率は72.3%となった。この結果から見て分かるとおり、LAサーバを利用することにより、ページの移動先発見数が21%増加した。

20

#### 【0049】

一方で、LAサーバを利用しても移動先の発見ができなかったものが13個存在した。そこで、実験において発見されたWebページ移動に伴うリンク切れ47個についてLAサーバが出力するログなどからLAサーバの探索結果についての分析を行った(図13)。

#### 【0050】

その結果、LAサーバがリンクオーソリティと考えられるページを発見でき、正しく更新されているページは47個中12個(図13e)であった。残る35個については図13a~dのように分析した。以下に分類毎の原因の検証を行う。

30

#### 【0051】

分類a:この分類に当てはまるものは18個存在した。これは今回の実験で利用したリンクオーソリティ候補収集の手法が原因であると考えられる。本来、リンクオーソリティ候補群Vとしては、uに対してリンクを貼っている全てのページが収集されるべきである。しかし、今回実験で用いたリンクオーソリティ候補の収集手法ではすべてのページを収集していないためである。

#### 【0052】

分類b:この分類に当てはまるものは14個存在した。これは上述のaと同様の理由も考えられる。しかし、全てのWebページにリンクオーソリティが存在するわけではないとも考えられる。つまり、リンクオーソリティ候補の収集手法の追加や改良を行ってもこのパターンに該当するものを限りなく減少させることは不可能であると考えられる。

40

#### 【0053】

分類c:実験ではこのパターンは存在していない。本実験の結果より、前述のリンクオーソリティ候補ランキングのヒューリスティクスが適切だったと考えることができる。この分類に当てはまるものが増加するようであれば、リンクオーソリティ候補のランキングヒューリスチクスの見直しが必要である。

#### 【0054】

分類d:この分類に当てはまるものは3個存在した。これは、Webページの移動が発生してから、リンクオーソリティが持つリンクを参照するまでの時間が短かったためと考

50

えられる。正しいリンクオーソリティは、Webページの移動に伴いリンクを正しく修正すると考えられるが、それがページ移動が起こった直後に行われるとは限らない。この分類に当てはまるものは、ページ移動後、ある程度時間が経過した後にリンクオーソリティが有効に機能するのではないかと考えられる。

【0055】

本実験では、特に分類 a と分類 b とに該当するものが多く存在した。この原因は、本実験で利用した LAサーバのリンクオーソリティ候補 V の収集手法に偏りがあったものと考えられる。リンクオーソリティ候補 v を効率的に収集するために、候補収集の手法を次に示す手順で行うことが好ましい。

【0056】

手順 1：全ての監視対象の URL に対して Alexa Web Information Service API (商標) を用いて、u へのリンクを持つページを収集。

【0057】

手順 2：全ての監視対象の URL に対して Google Web Service API (商標) を用いて、u へのリンクを持つページを収集。

【0058】

手順 3：全ての監視対象の URL に対してクローラを用いて u へのリンクを持つページがある可能性が高いと考えられる場所からの収集。

【0059】

上記実験結果では、手順 1 を用いた収集しか行っていない。そのため、候補の収集が必ずしも十分なものとならなかったものと考えられる。

【0060】

また、LAサーバによってリンクオーソリティと考えられるページが発見されていた 12 個については、ページ移動先の発見に及ぼした影響について検討した (図 14)。その結果、「リンクオーソリティが発見されていなければ、正しい移動先を上位 3 位以内にランキングできなかった」というものが 2 個存在した。つまり、LAサーバを利用することで、これまでの LIMサーバだけでは発見できない移動先を発見することができた。その一方で、「リンクオーソリティは発見できたが、正しい移動先を上位 3 位以内評価できなかった」というものが 1 個存在した。しかし、これは、LIMサーバによる移動先候補の評価手法の問題であり、LAサーバの問題ではなかった。

【0061】

上記の実験結果と検討から分かるように本発明の方法及び装置の実施の形態によれば、従来と比べて、短時間の内に高い精度でリンクオーソリティを決定できることが分かる。

【図面の簡単な説明】

【0062】

【図 1】リンクオーソリティの考え方を説明するための図である。

【図 2】リンクオーソリティ決定装置の概要を説明するために用いる図である。

【図 3】自動修正システムの構成 (アーキテクチャ) を示す図である。

【図 4】LIMサーバの基本動作を説明するためのアルゴリズムを示す図である。

【図 5】主として前述の LAサーバによって構成される本発明のリンクオーソリティ決定装置の実施の形態の構成の一例を示すブロック図である。

【図 6】リンクオーソリティ決定装置をコンピュータを用いて実現する場合に用いるプログラムのアルゴリズムを示すフローチャートである。

【図 7】リンクオーソリティ決定装置の処理を説明するために用いる図である。

【図 8】属性の内容を示す図である。

【図 9】間接相互リンクを説明するための図である。

【図 10】ランキング付けパターンの表を示す図である。

【図 11】実験期間内に発生したリンク切れの数を示す図である。

【図 12】移動先の探索結果を示す図である。

【図 13】リンクオーソリティの探索結果の分類を示す図である。

10

20

30

40

50

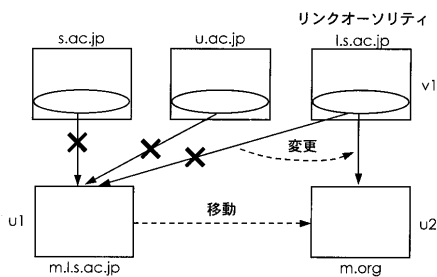
【図14】 移送先探索におけるリンクオーソリティを利用した影響を示す図である。

【符号の説明】

【0063】

- 1 リンクオーソリティ決定装置
- 3 Webページ収集手段
- 5 リンクオーソリティ候補決定手段
- 7 ランキング手段
- 9 属性記憶手段
- 11 候補決定手段
- 13 リンクオーソリティ決定手段
- 15 検索手段
- 17 ランキング手段
- 19 最終決定手段

【図1】



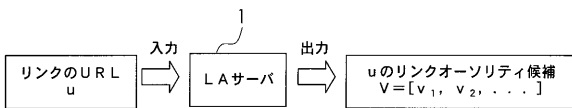
【図4】

```

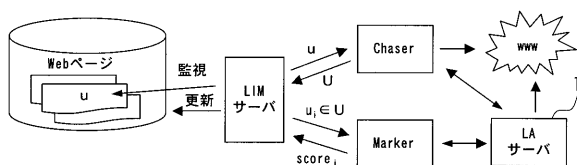
1.LIM-Server(u){
2. while(true){
3.  if(isBroken(u)){link is broken
4.  Set U=Chaser(u);
5.  for each ui in U{
6.  scorei=Maker(ui);
7.  }
8.  ListU=the result of scoring U by score ;
9.  }
10. interval();
11. }
12.}

```

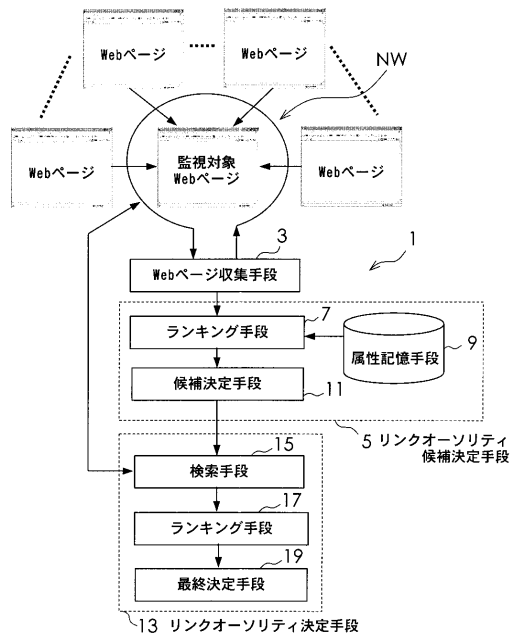
【図2】



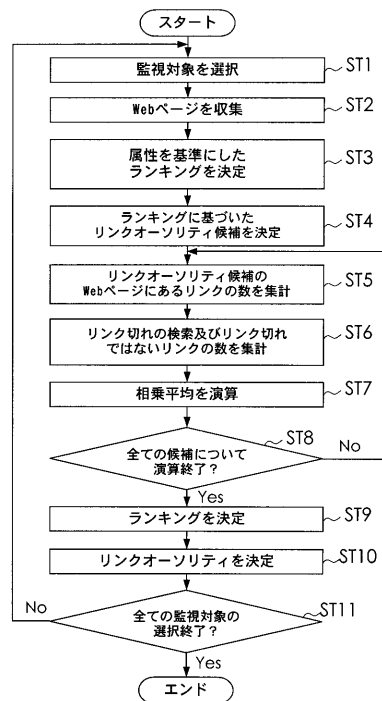
【図3】



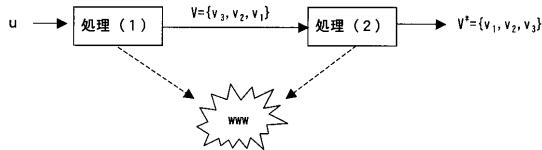
【 図 5 】



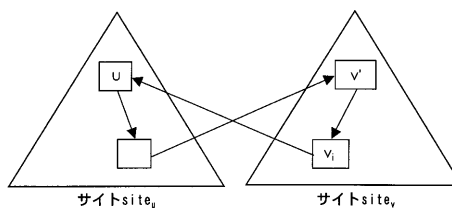
【 図 6 】



【 図 7 】



【 図 9 】



【 図 8 】

v <sub>i</sub> の属性	値	意味
同一同一	真偽	uと同一サイトの同一ディレクトリにv <sub>i</sub> が存在
同一上位	真偽	uと同一サイトの上位ディレクトリにv <sub>i</sub> が存在
同一下位	真偽	uと同一サイト且つ下位ディレクトリにv <sub>i</sub> が存在
同一その他	真偽	uと同一サイト且つその他のディレクトリにv <sub>i</sub> が存在
上位サイト	真偽	uが属するサイトをサブドメインとして含むサイトにv <sub>i</sub> が存在
外部サイト	真偽	uの上位サイト及び同一サイト以外のサイトにv <sub>i</sub> が存在
相互リンク	真偽	uとv <sub>i</sub> との間に直接的若しくは間接的な相互リンクが存在
index	真偽	v <sub>i</sub> のファイル名がデフォルトファイル名(典型的にはindex.html)である
#L	自然数	v <sub>i</sub> のページに含まれているリンクの数
B	[0, 1]	v <sub>i</sub> のページに含まれているリンクのうち、リンク切れではないものの割合

【 図 10 】

パターン	ランク	同一サイト				上位 サイト	下位 サイト	相互 リンク	index
		上位	同一	下位	その他				
a	1		●					●	-
b	2	●						●	-
c	2		●			●		●	-
d	2						●	●	-
e	3			●				●	-
f	4	●						●	-
g	5		●					●	-
h	6				●			●	-
i	7					●		●	-
j	7						●	●	-
k	8			●				-	-
l	8						●	-	-

【 図 11 】

監視対象ページ数	49750
リンク切れ総数	146
内、移動に伴うリンク切れ	47

## 【図 1 2】

	成功	失敗	成功率
LAサーバなし	28	19	59.6%
LAサーバあり	34	13	72.3%

## 【図 1 3】

分類	詳細	個数
a	$V = \phi$	18
b	$V \neq \phi$ だが、LAらしきページは含まれていない	14
c	VにLAと思われるページが含まれているが、 $V'$ の上位ではない	0
d	$V'$ の上位3位以内で発見されたが、まだリンクが更新されていない	3
e	$V'$ の上位5位以内で発見され、正しくリンクが変更された	12

## 【図 1 4】

効果	詳細	個数
あり	リンクオーソリティが発見されていなければ、正しい移動先を発見できなかった	4
	リンクオーソリティが発見されていなければ、正しい移動先を上位3位以内に評価できなかった	2
なし	リンクオーソリティは発見されたが、正しい移動先を上位3位以内に評価できなかった	1
	リンクオーソリティは発見されたが、移動先の発見に特に影響を与えなかった	5



---

フロントページの続き

(72)発明者 北川 博之  
茨城県つくば市天王台一丁目1番1 国立大学法人筑波大学内

(72)発明者 中溝 昌佳  
神奈川県川崎市川崎区貝塚2-7-9

Fターム(参考) 5B075 NK44  
5B082 HA08