

(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2006-23936

(P2006-23936A)

(43) 公開日 平成18年1月26日(2006.1.26)

(51) Int. Cl. F I テーマコード (参考)
G06F 17/30 (2006.01) G06F 17/30 340A 5B075
 G06F 17/30 180Z

審査請求 未請求 請求項の数 6 O L (全 10 頁)

<p>(21) 出願番号 特願2004-200717 (P2004-200717) (22) 出願日 平成16年7月7日(2004.7.7)</p>	<p>(71) 出願人 000233055 日立ソフトウェアエンジニアリング株式会社 神奈川県横浜市鶴見区末広町一丁目1番43 (74) 代理人 100091096 弁理士 平木 祐輔 (74) 代理人 100105463 弁理士 関谷 三男 (74) 代理人 100102576 弁理士 渡辺 敏章 (74) 代理人 100108394 弁理士 今村 健一</p>
---	---

最終頁に続く

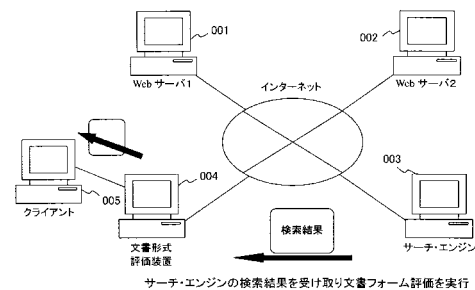
(54) 【発明の名称】 コンテンツ検索システム

(57) 【要約】

【課題】 ユーザーが指定した文書形式条件を満たすWebコンテンツを提供する。

【解決手段】 インターネットによって接続される第1 Webサーバ001、第2 Webサーバ002と、サーチ・エンジン003と、文書形式評価装置004と、クライアント端末005と、を有している。クライアント端末005は、例えばパーソナルコンピュータであり、ネットワークを介して第1又は第2 Webサーバ001、002上のWebコンテンツを閲覧する。クライアント端末005は、検索する語句及び文書形式に関する検索条件を、本実施の形態による文書形式評価装置004に送信する。文書形式評価装置004は、検索対象語句をサーチ・エンジン003に送信する。サーチ・エンジンは検索対象語句に基づいてコンテンツの検索を実行し、検索対象語句を含むWebコンテンツの一覧を文書形式評価装置004に返す。この一覧に含まれるWebコンテンツに対して文書形式を評価し、ユーザーが指定した条件に適合するコンテンツのみをユーザーへ返す。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

Webコンテンツを検索するコンテンツ検索システムであって、
コンテンツ検索の際にコンテンツの文書形式に基づいてフィルタリングを行う文書形式評価装置を有することを特徴とするコンテンツ検索システム。

【請求項 2】

クライアント端末及びコンテンツを提供するWebサーバとネットワークを介して接続され、Webコンテンツを検索するコンテンツ検索システムであって、
コンテンツの検索エンジンを備える検索装置と、

該コンテンツの検索結果として得られたコンテンツに関して文書形式条件に基づいて文書形式の評価を行う文書形式評価装置であって前記文書形式条件に合致したコンテンツを前記クライアント端末に提供する文書形式評価装置と
を有するコンテンツ検索システム。

10

【請求項 3】

クライアント端末及びコンテンツを提供するWebサーバとネットワークを介して接続され、Webコンテンツを検索するコンテンツ検索システムであって、

コンテンツの検索エンジンと該検索エンジンによるコンテンツの検索結果として得られたコンテンツに関して文書形式条件に基づいて文書形式の評価を行う文書形式評価部とを備える検索装置であって前記文書形式条件に合致したコンテンツを前記クライアント端末に提供する検索装置を、有するコンテンツ検索システム。

20

【請求項 4】

クライアント端末及びコンテンツを提供するWebサーバとネットワークを介して接続され、

文書形式条件に基づいて文書形式の評価を行う文書形式評価プログラムを保持する文書形式評価装置と、該文書形式評価装置における評価結果を蓄積するデータベースと、を備え、

前記文書形式評価プログラムが、前記ネットワークを介して複数の前記Webサーバ間を巡回しWebコンテンツの文書形式について評価を実行することを特徴とするコンテンツ検索システム。

【請求項 5】

前記評価結果を、前記検索を実行する前に前記評価を実行して前記データベースに保存しておくことを特徴とする請求項 4 に記載のコンテンツ検索システム。

30

【請求項 6】

請求項 1 から 5 までのいずれか 1 項に記載のコンテンツ検索システムであって、前記文書検索条件をコンテンツの検索キーに変換する変換部を備えるコンテンツ検索システム。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、コンテンツ検索システムに関し、特に、ユーザーが指定した文書形式条件を満たすWebコンテンツを検索するシステムに関する。

40

【背景技術】**【0002】**

インターネットによって接続されたWebサーバにおいては、非常に多くの情報を検索可能である反面、多くの情報の中から本当に閲覧したい情報を選択するのが難しいという問題がある。そこで、各種のフィルタ技術を用いて、不要な情報を検索対象から除外することがなされている。

【0003】

従来から、ユーザーがWebコンテンツを検索する場合に、使用するフィルタ項目としては、以下の項目が用いられていた。1) 興味対象となる語句、2) Webコンテンツのドメインで認識される組織、3) ファイルタイプ、4) ページ全体又はタイトルやリンク等の

50

ページ内容の一部が最終更新された時期等である。

【0004】

また、弱視や色盲等の視覚障害を持つユーザーや、キーボード及びマウス操作が困難なユーザーが、各々にとってのアクセシビリティ（利用容易性）を判断する基準項目をフィルタとして使用し、基準を満たすWebコンテンツを検索する方法が提案されている（特許文献1参照）。この特許文献1に記載の「利用容易性」とは、「障害や特別の必要性を持つユーザーが個別に定義するもの又はWWWコンソーシアムが公表するアクセシビリティ（利用容易性）・ガイドラインにおいて体系化されている定義」を指す。例えば、1）全ての画像について代替のテキストを提供する。2）前景色及び背景色は相互に十分なコントラストを有する。3）もし画像が、その代替テキストの情報以外に重要な情報を伝えるならば、拡張された記述を提供する。4）イベント・ハンドラは、マウスの使用を必要としないことを確認する等のフィルタ項目がある。本技術においては、視覚情報に依存するコンテンツや、マウス操作が不可能である場合に、ユーザーにとって利用し難いコンテンツを、ユーザーに提示する検索対象から除外するか、或いは、利用容易ではないことを示すインジケータとともに検索結果に表示することが可能である。

10

【0005】

【特許文献1】特開2002-334034号公報

【発明の開示】

【発明が解決しようとする課題】

【0006】

実際にWeb上で閲覧可能なコンテンツには画像などを主体としたコンテンツも多くなっているが、實際上、文書形式で記載されているコンテンツも非常に多いことは周知の通りである。さらに、一般のユーザーが閲覧する際に、画像データなどに関しては、サイズと容量とのみを変数となるため特に問題とならないが、文書データに関しては、種々の表示形式が存在し、ユーザーにとって見やすい文書形式ではない場合が多いのも事実である。

20

【0007】

上記特許文献1に記載の技術によってコンテンツが利用容易になる場合であっても、ユーザーが好む文書形式で提供されない場合には、閲覧を止めてしまうこともある。

【0008】

本発明は、コンテンツをユーザーの希望する文書形式で提供する技術に関する。

30

【課題を解決するための手段】

【0009】

本発明は、ユーザーが指定した文書形式条件を満たすコンテンツをユーザーへ提供する方法、又はプログラムによって構成されるWeb検索システムである。クライアント端末及びコンテンツを提供するWebサーバとはネットワークを介して接続可能な状態にあることが前提となる。本Web文書形式検索システムは、サーチ・エンジンが実行して得た検索結果を受け取り、検索結果に含まれるコンテンツが、ユーザーの指定した文書形式条件を満たすかどうかを評価する。条件を満たしていないと判断したコンテンツを、サーチ・エンジン検索結果から除外し、残りのコンテンツを結果セットとしてユーザーに返す。

40

【発明の効果】

【0010】

本発明により、ユーザーは、検索システムを使用して得られる莫大な量の検索結果から、好まない文書形式を持つコンテンツを除外した結果セットを得ることが可能となる。つまり、本発明によれば、ユーザーが検索結果として受け取るコンテンツの量を少なくすることで、結果としてネットワーク転送量を削減し、検索効率を向上させることが可能である。

【発明を実施するための最良の形態】

【0011】

本発明の実施の形態について説明する前に、発明者の行った考察について説明する。ユ

50

ユーザーが好まない文書形式を有するWebコンテンツを検索結果から除外しユーザーに表示する結果セットには含めないようにすることにより、より少数まで絞られた検索結果に基づいて、ユーザーが見やすいコンテンツに関する情報の提供を受けられることになる。このようなシステムを構築することにより、従来の方法による検索結果よりも、検索結果としてユーザーへ提供されるコンテンツ数及びコンテンツ量が削減される。従って、ネットワークの転送量が軽減されるという。また、ユーザーは「一見して閲覧を好まず閲覧を止める」ようなWebのページに、立ち入る機会自体が少なくなり検索効率を一層向上させることができる。

【0012】

ユーザーが好まない文書形式であると判断する項目の例を以下に挙げる。

10

1) 文字のサイズが大き過ぎる、または、小さ過ぎる。2) 文字の量に対して空白行の割合が少ない、または、行が一行又は数行おきに記述されていて空白が多過ぎる。3) 漢字の割合又はひらがなの割合が適正値から大きく外れている。或いは、カタカナが多過ぎたり、アルファベットが多過ぎたりする。4) 1ページの内容量が多く縦方向のスクロールを何度もする必要がある。5) ファイルサイズが大きく、例えば5秒以内に内容が表示されない。6) Webコンテンツの閲覧に使用している表示画面のサイズが、Webコンテンツの分量に比べて小さ過ぎるため、横方向のスクロールが必要である。7) ページの左右におけるマージンが少ない。8) ポップアップが多い。9) 背景に画像を使用している。10) Webコンテンツの閲覧に使用しているブラウザソフトウェアが、サポートしていない、または正常に表示することができない機能、ファイルのタイプを使用している。

20

【0013】

以上の1)から10)までの項目は、異なるユーザーの場合はむしろであるが、同じユーザーであっても、検索目的によって変化すると考える。例えば、語句の説明を得ることを検索目的とする場合は、簡潔な文章を検索結果とするために、フィルタ条件は多くなると考えられるが、語句に関する経験・感想等を趣味として閲覧する場合は、逆にフィルタ条件は少なくなることが好ましいと考えられる。また、検索に使用するパーソナルコンピュータのスペック(ネットワーク転送速度や画面サイズ、ブラウザソフトウェアの性能)によっても、上記の項目は変化すると考えられる。従って、フィルタ条件として使用する文書形式条件は、上記に挙げたような項目毎に設定できるようにしておき、且つ、簡易に設定変更できるようにすることが望ましい。

30

【0014】

以下、本発明の一実施の形態によるWeb文書形式検索システムについて、図面を参照しつつ説明を行う。

【0015】

図1は、本発明の一実施の形態によるWeb文書形式検索システムの概略構成例を示す図である。図1に示すWeb文書形式検索システムは、文書形式評価装置004がサーチ・エンジンとは独立に稼動する構成を示す図である。図1に示すように、本実施の形態によるWeb文書形式検索システムは、インターネットによって接続される第1Webサーバ001、第2Webサーバ002と、サーチエンジン003と、文書形式評価装置004と、クライアント端末005と、を有している。クライアント端末005は、例えばパーソナルコンピュータであり、ネットワークを介して第1又は第2Webサーバ001、002上のWebコンテンツを閲覧する。

40

【0016】

クライアント端末005は、検索する語句及び文書形式に関する検索条件を、本実施の形態による文書形式評価装置004に送信する。文書形式評価装置004は、検索対象語句をサーチ・エンジン003に送信する。サーチ・エンジンは検索対象語句に基づいてコンテンツの検索を実行し、検索対象語句を含むWebコンテンツの一覧を文書形式評価装置004に返す。この一覧に含まれるWebコンテンツに対して文書形式を評価し、ユーザーが指定した条件に適合するコンテンツのみをユーザーへ返す。

【0017】

50

尚、文書形式評価装置 004 に対して 1 又はそれ以上の端末が接続されており、文書形式評価装置 004 をプロキシ・サーバとして機能させる場合もある。

【0018】

次に、本発明の第 2 の実施の形態による Web 文書形式検索システムについて図面を参照しつつ説明を行う。図 2 は、本実施の形態によるシステムの構成例を示す図である。図 2 に示すように、本実施の形態によるシステムは、第 1 から第 3 までの Web サーバ 006 ~ 008 と、文書形式評価プログラムを搭載したサーチ・エンジン 009 と、クライアント端末 010 と、を有している。図 2 が図 1 と異なる点は、文書形式評価機能がサーチ・エンジンの一機能として稼動する点である。

【0019】

サーチ・エンジン 009 に文書形式評価装置を搭載しているため、図 1 に示すシステムと比較して、サーチ・エンジンと文書形式評価装置間のネットワーク転送が不要となる。従って、図 1 に示すシステムに比べて、ネットワーク転送量が少なくなり、ユーザーの検索時間を短縮することが可能となる。

【0020】

図 3 は、本発明の第 1 又は第 2 の実施の形態によるシステムにおいて、クライアント端末からの検索要求を受けて、クライアント端末に検索結果を返すまでの処理の流れを示すフローチャート図である。以下においては、図 1 の第 1 の実施の形態によるシステムを例に処理の流れを説明する。

【0021】

図 3 に示すように、処理を開始すると、まずステップ 101 において、Web コンテンツを検索しようとするクライアントは、クライアント端末を用いて Web 文書形式検索システムへ、A) 検索対象となる語句、及び、B) 検索条件とする文書形式(文書フォーム)に関する情報を送信する。ステップ 102 において、外部サーチ・エンジンが、A) 検索対象となる語句を含む Web コンテンツを検索し、ステップ 103 において、その検索結果を Web 文書形式検索システムの評価装置に渡す。

【0022】

ステップ 104 において、Web 文書形式検索システムは、検索結果であるコンテンツの文書形式を、検索条件として挙げられている項目について評価する。ステップ 105 において評価結果とユーザーが指定した条件とを比較し、B) 検索条件とする文書形式を満たす場合にはステップ 107 に、満たさない場合にはステップ 106 に進む。ステップ 106 においてクライアント要求に該当しないコンテンツをサーチ・エンジン検索結果から除外し、ステップ 107 においてクライアント要求に該当するコンテンツのみをクライアントへ送信する。本実施の形態による Web 文書形式検索システムは、サーチ・エンジンの一部であっても、独立したソフトウェアであっても良い。

【0023】

以下に、検索対象外とする文書形式条件候補となる各項目の例を挙げる。以下に挙げる具体的な数値については、クライアントがカスタマイズできる。

1) 文字のサイズが大き過ぎる(20px(ピクセル)以上、又は size 5 以上の文字が本文全体の 20% 以上。)

2) 文字のサイズが小さ過ぎる(10px 以下、又は size 2 以下の文字が本文全体の 40% 以上。)

3) 文字が多過ぎる(段落内の文字数が 300 文字以上の段落に含まれる文字数が本文全体の 50% 以上。)

4) 行が 1 行間隔で空けられている(<p>タグ属性 line-height: 200% (行間がフォントの高さの 2 倍) 以上の段落に含まれる文字数が本文全体の 30% 以上。)

5) 漢字が多過ぎる(本文全体の 60% 以上。)

6) カタカナが多過ぎる(本文全体の 60% 以上。)

7) アルファベットが多過ぎる(本文全体の 60% 以上。)

8) 1 ページの容量が多く縦方向のスクロールが数十回必要になる(ファイルサイズが

10

20

30

40

50

100KBを超えている。)。

9) ファイルサイズが大きく、5秒以内に内容が表示されない(ファイルサイズが、Webコンテンツ閲覧に使用しているシステムのネットワーク転送速度において、5秒で伝送できる量を超えている。)。

10) Webコンテンツの閲覧に使用しているパーソナルコンピュータの画面サイズが、Webコンテンツに比べて小さい為、横方向のスクロールが必要である(テーブル及び画像の横幅が600ドット以上。)。

11) ページの左右にマージンが無い(bodyタグ属性margin)。

12) ポップアップが多い(2個以上。)。

13) 背景に画像を使用している(bodyタグ属性backgroundが指定されている。)。 10

14) 背景に原色が使用されている(例: 赤色 bodyタグ属性bgcolor=#ff0000)。

15) 背景が白色ではない(bodyタグ属性bgcolor=#ffffffではない。)。

16) Webコンテンツの閲覧に使用しているブラウザソフトウェアが、サポートしていない、又は、表示に不安がある機能やファイルのタイプを含む。

A.PNG (Portable Network Graphics) 形式の画像を含む。

B.スタイルシートを使用している。

C.フレームを使用している。

【0024】

クライアントは、検索を実行する前に、上記の項目のそれぞれについて、検索条件として使用するか否か、また、各項目の数値等を決定する。上記項目のうち、条件1)~7) 20
は、コンテンツ内のテキスト本文の見易さ、読み易さに関する条件を示す項目である。以下、条件4)を設定した場合を例にし、図5を参照して具体的判定手順を説明する。まず、対象となるコンテンツのHTMLファイルを取得する(ステップ501)。全文字数、見難い文字数にゼロを初期設定した上で(ステップ502)、HTMLファイル内のタグ(<>で囲まれた部分を、全部のタグが終了するまで順に処理する(S503、S504)。

【0025】

次に現れたタグが<p>タグで、かつline-height属性が200%以上(行間がフォントの高さの2倍以上)の場合は、見難い文字であるとして、<p>タグが終了するまで(</p>が現れるまで)に出現する本文(<>で囲まれていない文字)の文字数をカウントし、全文字数と見難い文字数のそれぞれをカウントアップする(S506)。それ以外の場合(<p>タグ以外または<p>タグであってもline-height属性が200%未満と解釈される場合)は、タグが終了するまで全文字数のみをカウントアップする(S507)。 30

【0026】

このようにして全部のタグの処理が終わった時点で、見難い文字数が全文字数に占める割合を計算し(S508)、30%以上の場合は見難いコンテンツとして当該コンテンツを検索対象から除外し(S510)、30%未満の場合は当該コンテンツを検索対象に含める(S509)。既に述べた通り、この200%、30%等の数値は、クライアント端末からクライアントが入力等できるカスタマイズ可能な値である。

【0027】

条件4)以外の他の条件についても同様に、それぞれの条件に合致した「見難い文字数」をカウントし、本文の全文字数中に占める割合が一定以上の場合には見難いコンテンツであると判定して、当該コンテンツを検索対象から除外する。尚、条件3)の判定手順において、段落の開始、終了の判断は、<P>タグで判断する方法の他に、例えば
タグ(改行タグ)が出現したら、あるいは本文を挟まずに2つ以上続いたら段落が切り変わったというような判断をしてもよい。また、上記項目のうち、例えば条件8)~16)は、クライアント端末、ネットワークの性能等に依存する見難さなど、文字自体の読み難さ以外の要素に依存する条件を示す項目である。 40

【0028】

例えば、ネットワーク転送速度36kbpsのモデムを介してネットワークに接続した 50

B 5 サイズのノート型パソコンを使用し、興味対象となる語句の意味を検索する目的で、Webコンテンツを閲覧しているユーザーがいるとする。ユーザーは、意味の説明として画像が検索結果に含まれることを期待しているが、使用しているブラウザはPNG形式の画像表示をサポートしていないとする。このような場合には、まず、上記項目のうち以下の文書形式コンテンツを検索対象外に指定することで、ネットワーク転送量の負荷を抑えることが可能である。

9) ページの容量が多く縦方向のスクロールが数十回必要になる(ファイルサイズが100KBを超えている)。

10) ポップアップが多い(例えば2個以上)。

14) 背景に画像を使用している。

10

【0029】

また、以下の項目を検索対象外に指定することで、表示できない画像を検索結果に含めないようにすることが可能である。

17) A.PNG (Portable Network Graphics) 形式の画像を含む。

【0030】

以下の項目を検索対象外に指定することで、B 5 サイズのノートパソコンを使用しても横方向にスクロールする必要が無く、見易さを感じるコンテンツのみ、検索結果として受け取ることができる。

11) テーブル及び画像の横幅が600ドット以上。

16) 背景が白色ではない(bodyタグ属性bgcolor=#ffffffではない)。

20

Web文書形式検索システムは、プロキシ・サーバ上ではなくとも、クライアントパーソナルコンピュータ上のソフトウェアとして実行されても良い。プロキシ・サーバ上のソフトウェアとしてWeb文書形式検索システムが稼動する場合、Webコンテンツの文書形式について評価を実行するタイミングとしては、前述のように「サーチ・エンジンからの検索結果を受け取った後」の他に、定期的に評価プログラムがWebサーバ上の全Webコンテンツを対象として評価を実行し、上記に挙げた全項目についての評価結果をデータベース化しておくという方法でも良い。この方法について第3の実施の形態として図4を参照しつつ説明する。

【0031】

図4に示すように、本実施の形態によるWeb文書形式検索システムは、インターネットと、第1のWebサーバ011と、第2のWebサーバ012と、第3のWebサーバ013と、文書形式評価装置014と、データベース016と、を有する。

30

【0032】

本実施の形態によるWeb文書形式検索システムにおいては、文書形式評価プログラム015は、定期的に第1のWebサーバ011、第2のWebサーバ012、第3のWebサーバ013の間を、これらを接続するネットワークを介して巡回し、Webコンテンツの文書形式について評価を実行する。評価結果を、文書形式評価サーバ014が所有するデータベース016に保存する。ユーザーが検索を実行した後ではなく、事前に評価を実行しておくことにより、Web文書形式検索システムにおける検索速度が向上する。

【0033】

以上、本発明について各実施の形態について説明したが、本発明はこれらに限定されるものではなく、種々の変形が可能である。例えば、ユーザーの要望に応じて文書形式を変換する変換部を有し、文書形式を変換した後にユーザー端末に提供するようにすることも可能である。或いは、文書形式に関してユーザーの要望(文書形式条件)と実際のコンテンツにおける文書形式との偏差を演算し、ある程度までの偏差を許容するようにすることも可能である。により、ユーザーは、従来の検索システムを使用して得られる莫大な量の検索結果から、好まない文書形式を持つコンテンツを除外した結果セットを得ることが可能となる。つまり、本発明は、ユーザーが検索結果として受け取るコンテンツの量を減少させ、結果としてネットワーク転送量を削減し、ユーザーの検索効率を向上させることが可能である。

40

50

【図面の簡単な説明】

【0034】

【図1】本発明の第1の実施の形態によるWebコンテンツ検索システムWebコンテンツ検索システムがサーチ・エンジンと独立して稼動する場合の構成例を示す図である。

【図2】本発明の第2の実施の形態によるWebコンテンツ検索システムの構成例を示す図であり、Webコンテンツ検索システムがサーチ・エンジンの一機能として稼動する場合の例を示す図である。

【図3】本発明の実施の形態によるWebコンテンツ検索システムにおける処理の流れを示すフローチャート図である。

【図4】本発明の第3の実施の形態によるWebコンテンツ検索システムであって、評価プログラムがWebサーバ上の全Webコンテンツを対象として、定期的に文書形式についての評価を実行し、その評価結果を記録しておくシステムの概念である。

【図5】行が1行間隔で空けられている(<p>タグ属性 line-height: 200% (行間がフォントの高さの2倍以上)の段落に含まれる文字数が本文全体の30%以上。)という条件に基づくコンテンツ検索処理の流れを示すフローチャート図である。

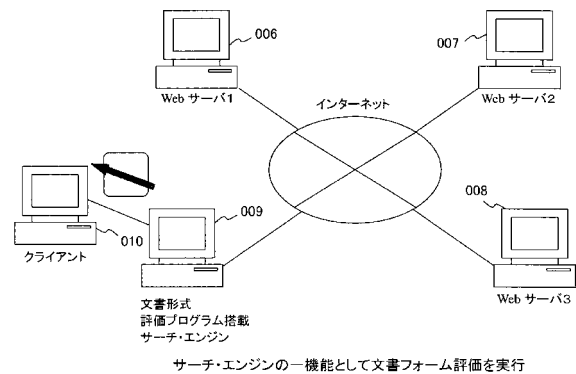
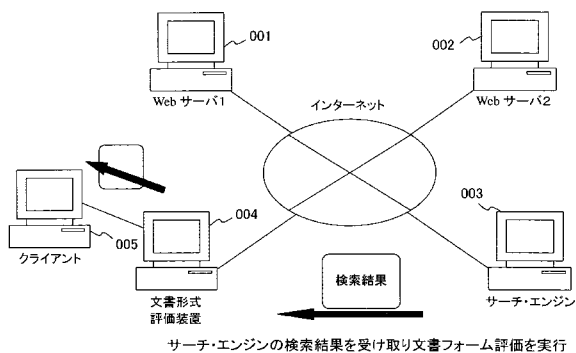
【符号の説明】

【0035】

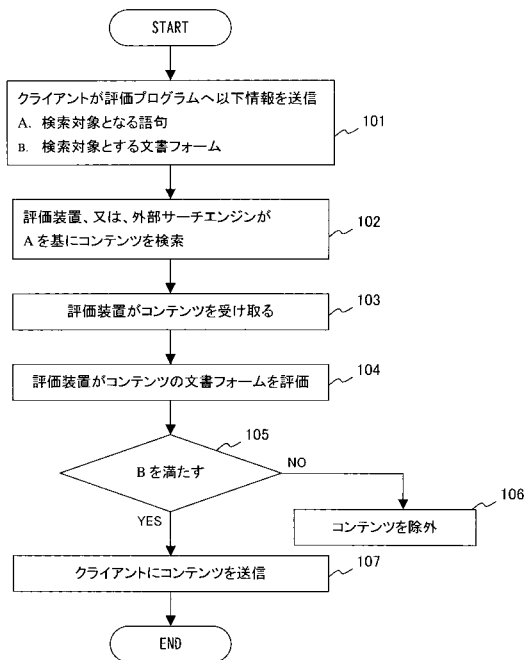
001/002...Webサーバ、003...サーチ・エンジン、004...文書フォーム評価装置、005...クライアントパーソナルコンピュータ、006/007/008...Webサーバ、009...文書形式評価プログラム搭載サーチ・エンジン、010...クライアントパーソナルコンピュータ、011/012/013...Webサーバ、014...文書形式評価装置、015...Webコンテンツの文書形式に関する評価実行プログラム、016...015による評価結果を格納するデータベース。

【図1】

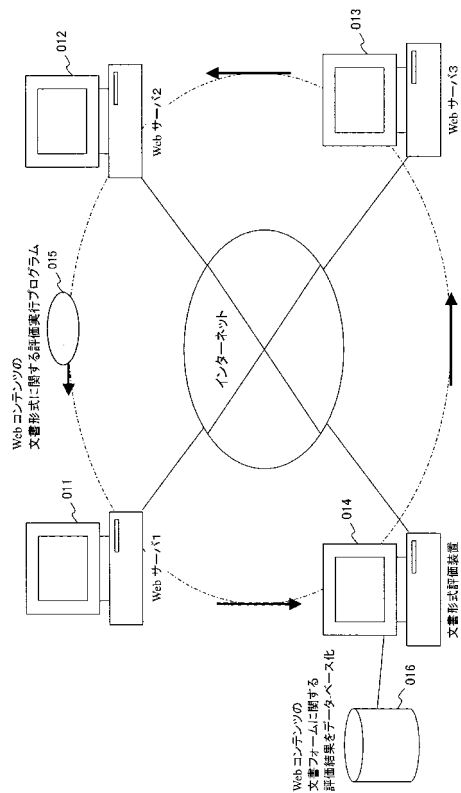
【図2】



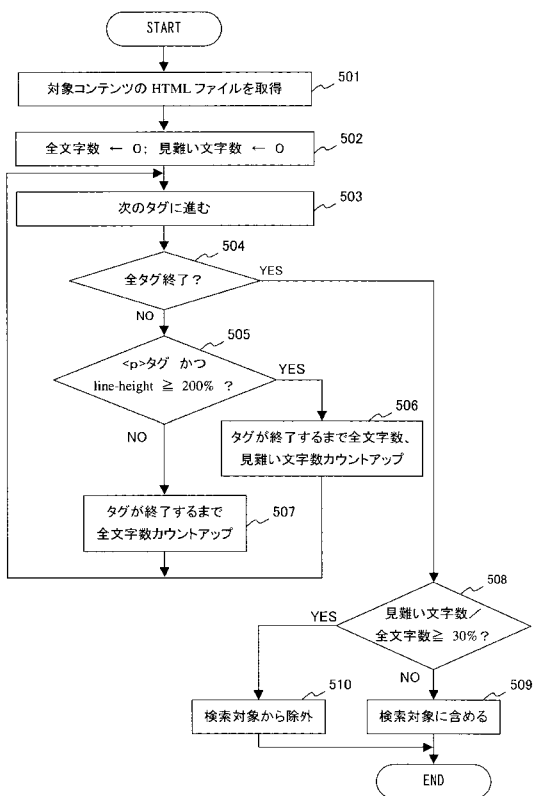
【 図 3 】



【 図 4 】



【 図 5 】



フロントページの続き

(72)発明者 松下 智子

東京都品川区東品川4丁目1番7号 日立ソフトウェアエンジニアリング株式会社内

Fターム(参考) 5B075 KK02 KK07 ND03 ND16 NR02 PQ02 PQ42 PR08