

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5462591号
(P5462591)

(45) 発行日 平成26年4月2日(2014.4.2)

(24) 登録日 平成26年1月24日(2014.1.24)

(51) Int.Cl.		F I			
G06Q	30/02	(2012.01)	G06Q	30/02	150
G06F	13/00	(2006.01)	G06F	13/00	540P
G06F	17/30	(2006.01)	G06F	17/30	210A

請求項の数 11 (全 25 頁)

(21) 出願番号	特願2009-250646 (P2009-250646)	(73) 特許権者	399037405 楽天株式会社 東京都品川区東品川四丁目12番3号
(22) 出願日	平成21年10月30日(2009.10.30)	(74) 代理人	110000958 特許業務法人 インテクト国際特許事務所
(65) 公開番号	特開2011-96078 (P2011-96078A)	(72) 発明者	岩淵 志学 東京都品川区東品川四丁目12番3号 楽天株式会社内
(43) 公開日	平成23年5月12日(2011.5.12)		
審査請求日	平成23年11月7日(2011.11.7)		
審判番号	不服2012-21010 (P2012-21010/J1)		
審判請求日	平成24年10月24日(2012.10.24)		
早期審査対象出願		合議体	
		審判長	手島 聖治
		審判官	須田 勝巳
		審判官	清田 健一

最終頁に続く

(54) 【発明の名称】 特有コンテンツ判定装置、特有コンテンツ判定方法、特有コンテンツ判定プログラム及び関連コンテンツ挿入装置

(57) 【特許請求の範囲】

【請求項1】

所定のサイトに含まれる複数のWebページのうち、指定されたWebページを構成しているコンテンツを抽出する抽出手段と、

前記複数のWebページのうち、前記指定されたWebページを構成している各コンテンツが他のWebページで用いられる頻度をカウントする計算手段と、

前記指定されたWebページを構成しているコンテンツのうち、他のWebページで用いられる頻度が所定値以下のコンテンツを当該指定されたWebページに特有のコンテンツであると判定する判定手段と、

を備えることを特徴とする特有コンテンツ判定装置。

10

【請求項2】

請求項1に記載の特有コンテンツ判定装置において、

前記抽出手段は、1つ以上のコンテンツで構成されたコンテンツグループの単位で、Webページを構成しているコンテンツを抽出し、

前記計算手段は、前記指定されたWebページを構成しているコンテンツグループが他のWebページで用いられる頻度をカウントし、

前記判定手段は、前記指定されたWebページを構成しているコンテンツグループのうち、他のWebページで用いられる頻度が所定値以下のコンテンツグループを当該指定されたWebページに特有のコンテンツグループであると判定することを特徴とする特有コンテンツ判定装置。

20

【請求項3】

請求項2に記載の特有コンテンツ判定装置において、
前記抽出手段は、所定のマークアップ言語で記述され、Webページを構成するコンテンツを示すドキュメントデータに基づいて、コンテンツグループを抽出することを特徴とする特有コンテンツ判定装置。

【請求項4】

請求項3に記載の特有コンテンツ判定装置において、
前記抽出手段は、前記コンテンツを示すドキュメントデータにおいて予め定められたタグに基づいてコンテンツグループを定めることを特徴とする特有コンテンツ判定装置。

【請求項5】

請求項1乃至4の何れか1項に記載の特有コンテンツ判定装置において、
前記抽出手段は、投稿された記事が掲載されるWebページから前記記事に対して投稿されたコメントを抽出し、

前記抽出された各コメントを、コメントが示す内容別に分類する分類手段と、
出現頻度の閾値を設定する設定手段であり、前記コメントが分類された前記内容の数が
多いほど前記閾値を小さくする設定手段と、

を更に備え、
前記計算手段は、前記コメントが分類された各前記内容のWebページにおける出現頻度を計算し、

前記判定手段は、前記計算手段により計算された出現頻度が前記設定された閾値以下である前記内容を前記Webページに特有の内容であると判定することを特徴とする特有コンテンツ判定装置。

【請求項6】

所定のサイトに含まれる複数のWebページのうち、指定されたWebページを構成しているコンテンツを抽出する抽出工程と、

前記複数のWebページのうち、前記指定されたWebページを構成している各コンテンツが他のWebページで用いられる頻度をカウントする計算工程と、

前記指定されたWebページを構成しているコンテンツのうち、他のWebページで用いられる頻度が所定値以下のコンテンツを当該指定されたWebページに特有のコンテンツであると判定する判定工程と、

を有することを特徴とする特有コンテンツ判定方法。

【請求項7】

請求項6に記載の特有コンテンツ判定方法において、
前記抽出工程は、投稿された記事が掲載されるWebページから前記記事に対して投稿されたコメントを抽出し、

前記抽出された各コメントを、コメントが示す内容別に分類する分類工程と、
出現頻度の閾値を設定する設定工程であり、前記コメントが分類された前記内容の数が
多いほど前記閾値を小さくする設定工程と、

を更に含み、
前記計算工程は、前記コメントが分類された各前記内容のWebページにおける出現頻度を計算し、

前記判定工程は、前記計算工程により計算された出現頻度が前記設定された閾値以下である前記内容を前記Webページに特有の内容であると判定することを特徴とする特有コンテンツ判定方法。

【請求項8】

コンピュータを、
所定のサイトに含まれる複数のWebページのうち、指定されたWebページを構成しているコンテンツを抽出する抽出手段、

前記複数のWebページのうち、前記指定されたWebページを構成している各コンテンツが他のWebページで用いられる頻度をカウントする計算手段、及び、

10

20

30

40

50

前記指定されたW e b ページを構成しているコンテンツのうち、他のW e b ページで用いられる頻度が所定値以下のコンテンツを当該指定されたW e b ページに特有のコンテンツであると判定する判定手段、

として機能させることを特徴とする特有コンテンツ判定プログラム。

【請求項 9】

請求項 8 に記載の特有コンテンツ判定プログラムにおいて、

前記抽出手段は、投稿された記事が掲載されるW e b ページから前記記事に対して投稿されたコメントを抽出し、

前記コンピュータを、

前記抽出された各コメントを、コメントが示す内容別に分類する分類手段、及び、

出現頻度の閾値を設定する設定手段であり、前記コメントが分類された前記内容の数が多いほど前記閾値を小さくする設定手段、

として更に機能させ、

前記計算手段は、前記コメントが分類された各前記内容のW e b ページにおける出現頻度を計算し、

前記判定手段は、前記計算手段により計算された出現頻度が前記設定された閾値以下である前記内容を前記W e b ページに特有の内容であると判定することを特徴とする特有コンテンツ判定プログラム。

【請求項 10】

請求項 1 乃至 5 の何れか 1 項に記載の特有コンテンツ判定装置と、

前記特有コンテンツ判定装置により特有のコンテンツであると判定されたコンテンツに関連する関連コンテンツを、前記指定されたW e b ページに挿入する挿入手段と、

を備えることを特徴とする関連コンテンツ挿入装置。

【請求項 11】

請求項 10 に記載の関連コンテンツ挿入装置において、

前記特有コンテンツ判定装置は、前記指定されたW e b ページを構成しているコンテンツとして、投稿された記事のテキストデータが含まれている場合に、当該テキストデータを当該W e b ページに特有のコンテンツであると判定し、

前記特有コンテンツ判定装置により特有のコンテンツであると判定された記事のテキストデータから前記指定されたW e b ページの特徴語を抽出する特徴語抽出手段と、

それぞれ語に関連付けて記憶手段に記憶された複数のコンテンツの中から、前記抽出された特徴語に関連するコンテンツを前記関連コンテンツとして選択する選択手段と、

を更に備え、

前記挿入手段は、前記選択された関連コンテンツを、前記指定されたW e b ページに挿入することを特徴とする関連コンテンツ挿入装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、W e b ページを構成するコンテンツを抽出する技術分野に関する。

【背景技術】

【0002】

従来、W e b サイト上に公開されているW e b ページを構成している素材たるコンテンツを取得し、取得したコンテンツに基づいて新たなコンテンツを生成する技術が知られている。例えば、非特許文献 1 には、ユーザにより画像データのU R L が指定されると、当該U R L に対応する画像データをW e b 上から取得し、取得した画像データに基づいてバナーを自動作成する技術が開示されている。

【先行技術文献】

【非特許文献】

【0003】

【非特許文献 1】 “バナー自動作成”、[online]、[平成21年10月21日検索]、インターネ

10

20

30

40

50

ット<U R L : http://hyperbannermaker.com/>

【発明の概要】

【発明が解決しようとする課題】

【0004】

Webサイトを構成する各Webページには、そのWebサイトの目的に沿った内容が掲載される。そのため、Webサイトを構成する各Webページの内容は、基本的には互いに関連性を有しているのであるが、夫々何らかの特徴を有してる場合がある。そして、そのWebページの内容を決める要因が、Webページを構成しているコンテンツ（例えば、テキストデータ、画像データ等）の内容である。従って、Webページを構成しているコンテンツの中で、そのWebページを特徴付けるコンテンツ、すなわち、そのWebページ特有のコンテンツが存在することがある。

10

【0005】

非特許文献1に記載の技術は、Webページに特有のコンテンツを抽出するものであるが、自動的に抽出するものではなく、ユーザが手作業でコンテンツを指定しなければならず、該Webページに特有のコンテンツを容易に抽出することはできない。そのため、どのコンテンツがWebページ特有のコンテンツであるかをユーザが判断することができない場合、又はユーザの嗜好によって好みのコンテンツが偏ってしまう場合など、Webページに特有のコンテンツを適格に抽出することができない。また、対象とするWebページのページ数が多いと、ユーザの作業が甚大となってしまう問題があった。

【0006】

20

また、例えば、HTML（HyperText Markup Language）文書のタグの記述に基づいて、画像だけ又はテキストだけといったように、特定種類のコンテンツを全て抽出することは可能である。しかしながら、抽出されたコンテンツの中には、Webページ特有のものではないありふれたコンテンツも含まれるので、抽出結果としての信頼性が低く、そのため抽出結果からユーザが特有のコンテンツを探さなければならなかった。

【0007】

本発明は以上の点に鑑みてなされたものであり、Webページを構成しているコンテンツの中からのそのWebページ特有のコンテンツを容易に抽出することができる特有コンテンツ判定装置、特有コンテンツ判定方法、特有コンテンツ判定プログラム等を提供することを目的とする。

30

【課題を解決するための手段】

【0008】

上記課題を解決するために、請求項1に記載の発明は、所定のサイトに含まれる複数のWebページのうち、指定されたWebページを構成しているコンテンツを抽出する抽出手段と、前記複数のWebページのうち、前記指定されたWebページを構成している各コンテンツが他のWebページで用いられる頻度をカウントする計算手段と、前記指定されたWebページを構成しているコンテンツのうち、他のWebページで用いられる頻度が所定値以下のコンテンツを当該指定されたWebページに特有のコンテンツであると判定する判定手段と、を備えることを特徴とする。

【0009】

40

この発明によれば、所定のサイトに含まれる複数のWebページのうち、指定されたWebページを構成している各コンテンツが複数のWebページで用いられる頻度が夫々カウントされる。他のWebページで用いられる頻度が小さいコンテンツであるほど、指定されたWebページ以外にはあまり出現しないコンテンツである。そこで、他のWebページで用いられる頻度が所定値以下であるかを判定することで、当該条件を満たす全てのコンテンツが、指定されたWebページに特有のコンテンツであると特定される。よって、Webページに特有のコンテンツを容易に抽出することができる。

【0012】

請求項2に記載の発明は、請求項1に記載の特有コンテンツ判定装置において、前記抽出手段は、1つ以上のコンテンツで構成されたコンテンツグループの単位で、Webページ

50

ジを構成しているコンテンツを抽出し、前記計算手段は、前記指定されたWebページを構成しているコンテンツグループが他のWebページで用いられる頻度をカウントし、前記判定手段は、前記指定されたWebページを構成しているコンテンツグループのうち、他のWebページで用いられる頻度が所定値以下のコンテンツグループを当該指定されたWebページに特有のコンテンツグループであると判定することを特徴とする。

【0013】

この発明によれば、コンテンツグループの単位でWebページに特有のコンテンツが判断されるので、例えば、Webページ上において或るまとまりをもって表示されていたり、互いに関連性を有しているようなコンテンツをコンテンツグループとしたときに、Webページに特有のコンテンツとなるものを抽出することができる。

10

【0014】

請求項3に記載の発明は、請求項2に記載の特有コンテンツ判定装置において、前記抽出手段は、所定のマークアップ言語で記述され、Webページを構成するコンテンツを示すドキュメントデータに基づいて、コンテンツグループを抽出することを特徴とする。

【0015】

この発明によれば、Webページを構成するコンテンツを示すドキュメントデータに基づいてコンテンツグループが抽出されるので、適格にコンテンツグループを抽出することができる。

【0016】

請求項4に記載の発明は、請求項3に記載の特有コンテンツ判定装置において、前記抽出手段は、前記コンテンツを示すドキュメントデータにおいて予め定められたタグに基づいてコンテンツグループを定めることを特徴とする。

20

【0017】

この発明によれば、予め定められたタグに基づいてコンテンツグループが抽出されるので、Webページに特有のコンテンツと、特有ではないコンテンツとが夫々予め定められたタグでグループ化されている場合に、Webページに特有のコンテンツを判断する精度を上げることができる。

【0018】

請求項5に記載の発明は、請求項1乃至4の何れか1項に記載の特有コンテンツ判定装置において、前記抽出手段は、投稿された記事が掲載されるWebページから前記記事に対して投稿されたコメントを抽出し、前記抽出された各コメントを、コメントが示す内容別に分類する分類手段と、出現頻度の閾値を設定する設定手段であり、前記コメントが分類された前記内容の数が多いほど前記閾値を小さくする設定手段と、を更に備え、前記計算手段は、前記コメントが分類された各前記内容のWebページにおける出現頻度を計算し、前記判定手段は、前記計算手段により計算された出現頻度が前記設定された閾値以下である前記内容を前記Webページに特有の内容であると判定することを特徴とする。

30

請求項6に記載の発明は、所定のサイトに含まれる複数のWebページのうち、指定されたWebページを構成しているコンテンツを抽出する抽出工程と、前記複数のWebページのうち、前記指定されたWebページを構成している各コンテンツが他のWebページで用いられる頻度をカウントする計算工程と、前記指定されたWebページを構成しているコンテンツのうち、他のWebページで用いられる頻度が所定値以下のコンテンツを当該指定されたWebページに特有のコンテンツであると判定する判定工程と、を有することを特徴とする。

40

請求項7に記載の発明は請求項6に記載の特有コンテンツ判定方法において、前記抽出工程は、投稿された記事が掲載されるWebページから前記記事に対して投稿されたコメントを抽出し、前記抽出された各コメントを、コメントが示す内容別に分類する分類工程と、出現頻度の閾値を設定する設定工程であり、前記コメントが分類された前記内容の数が多いほど前記閾値を小さくする設定工程と、を更に含み、前記計算工程は、前記コメントが分類された各前記内容のWebページにおける出現頻度を計算し、前記判定工程は、前記計算工程により計算された出現頻度が前記設定された閾値以下である前記内容を前記

50

Web ページに特有の内容であると判定することを特徴とする。

【0019】

請求項 8 に記載の発明は、コンピュータを、所定のサイトに含まれる複数の Web ページのうち、指定された Web ページを構成しているコンテンツを抽出する抽出手段、前記複数の Web ページのうち、前記指定された Web ページを構成している各コンテンツが他の Web ページで用いられる頻度をカウントする計算手段、及び、前記指定された Web ページを構成しているコンテンツのうち、他の Web ページで用いられる頻度が所定値以下のコンテンツを当該指定された Web ページに特有のコンテンツであると判定する判定手段、として機能させることを特徴とする。

請求項 9 に記載の発明は、請求項 8 に記載の特有コンテンツ判定プログラムにおいて、前記抽出手段は、投稿された記事が掲載される Web ページから前記記事に対して投稿されたコメントを抽出し、前記コンピュータを、前記抽出された各コメントを、コメントが示す内容別に分類する分類手段、及び、出現頻度の閾値を設定する設定手段であり、前記コメントが分類された前記内容の数が多いほど前記閾値を小さくする設定手段、として更に機能させ、前記計算手段は、前記コメントが分類された各前記内容の Web ページにおける出現頻度を計算し、前記判定手段は、前記計算手段により計算された出現頻度が前記設定された閾値以下である前記内容を前記 Web ページに特有の内容であると判定することを特徴とする。

【0020】

請求項 11 に記載の発明は、請求項 1 乃至 6 の何れか 1 項に記載の特有コンテンツ判定装置と、前記特有コンテンツ判定装置により特有のコンテンツであると判定されたコンテンツに関連する関連コンテンツを、前記指定された Web ページに挿入する挿入手段と、を備えることを特徴とする。

【0021】

この発明によれば、特有のコンテンツであると判定されたコンテンツに関連するコンテンツが、指定された Web ページに挿入されるので、Web ページの特徴と関連する情報を当該 Web ページに追加することができる。

【0022】

請求項 12 に記載の発明は、請求項 11 に記載の関連コンテンツ挿入装置において、前記特有コンテンツ判定装置は、前記指定された Web ページを構成しているコンテンツとして、投稿された記事のテキストデータが含まれている場合に、当該テキストデータを当該 Web ページに特有のコンテンツであると判定し、前記特有コンテンツ判定装置により特有のコンテンツであると判定された記事のテキストデータから前記指定された Web ページの特徴語を抽出する特徴語抽出手段と、それぞれ語に関連付けて記憶手段に記憶された複数のコンテンツの中から、前記抽出された特徴語に関連するコンテンツを前記関連コンテンツとして選択する選択手段と、を更に備え、前記挿入手段は、前記選択された関連コンテンツを、前記指定された Web ページに挿入することを特徴とする。

【0023】

この発明によれば、各記事のテキストデータにその記事特有の内容が含まれているのであれば、特有コンテンツ判定装置により各記事のテキストデータを抽出することができる。これにより、Web ページに掲載されている記事の内容に関連する情報を当該 Web ページに追加することができる。

【発明の効果】

【0024】

本発明によれば、他の Web ページで用いられる頻度が小さいコンテンツであるほど、指定された Web ページ以外に前記複数の Web ページ上にあまり出現しないコンテンツであるので、他の Web ページで用いられる頻度が所定値以下であるかを判定することで、当該条件を満たす全てのコンテンツが、指定された Web ページに特有のコンテンツであると特定される。よって、Web ページに特有のコンテンツを容易に抽出することができる。

10

20

30

40

50

【図面の簡単な説明】

【0025】

【図1】一実施形態に係るブログシステムSの概要構成の一例を示す図である。

【図2】一実施形態に係るブログサーバ1の概要構成の一例を示すブロック図である。

【図3】ブロガーが指定されてからブログページに広告コンテンツが挿入されるまでの処理の概要を示す図である。

【図4】Webページの構成例を示す図である。

【図5】HTML文書から生成されたDOMツリーの一例を示す図である。

【図6】記憶部15に記憶されたコンテンツブロック対応情報の内容の一例を示す図である。

【図7】一実施形態に係るブログサーバ1のシステム制御部20の広告コンテンツ挿入処理における処理例を示すフローチャートである。

【図8】一実施形態に係るブログサーバ1のシステム制御部20の1ページ対応抽出処理における処理例を示すフローチャートである。

【図9】一実施形態に係るブログサーバ1のシステム制御部20のツリー探索処理における処理例を示すフローチャートである。

【図10】一実施形態に係るブログサーバ1のシステム制御部20の特有コンテンツブロック判定処理における処理例を示すフローチャートである。

【図11】一実施形態の変形例に係るブログサーバ1のシステム制御部20のブログ更新時処理における処理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】

【0026】

以下、図面を参照して本発明の実施形態について詳細に説明する。なお、以下に説明する実施の形態は、ブログサービスを提供するブログシステムにおいて、ブログページ送信するサーバ装置に対して本発明を適用した場合の実施形態である。

【0027】

[1. ブログシステムの構成及び機能概要]

先ず、本実施形態に係るブログシステムSの構成及び概要機能について、図1を用いて説明する。

【0028】

図1は、本実施形態に係るブログシステムSの概要構成の一例を示す図である。

【0029】

図1に示すように、ブログシステムSは、特有コンテンツ判定装置及び関連コンテンツ挿入装置の一例としてのブログサーバ1と、管理端末2と、複数のユーザ端末3と、を含んで構成されている。そして、ブログサーバ1と、各ユーザ端末3とは、ネットワークNWを介して、例えば、通信プロトコルにTCP/IP等を用いて相互にデータの送受信が可能になっている。なお、ネットワークNWは、例えば、インターネット、専用通信回線（例えば、CATV（Community Antenna Television）回線）、移動体通信網（基地局等を含む）、及びゲートウェイ等により構築されている。また、ブログサーバ1と管理端末2とは、LAN（Local Area Network）等のネットワークを介して接続されている。

【0030】

このような構成のブログシステムSにおいて、ブログサーバ1は、ユーザ端末3からのリクエストに応じて、ブログサービスサイトを構成するWebページを送信するWebサーバである。ユーザ端末3を利用するユーザがブログサービスサイトのユーザ登録をすると、そのユーザは、当該ブログサービスサイトにおいてユーザ自身のブログを運営することができるようになっていく。そして、ユーザ登録されたユーザ（ブロガー）は、ブログサービスサイトにアクセスして、自己のブログを更新（ブログ記事（ブログ1件毎の記録）を追加）することができるようになっていく。そのため、ブログサーバ1は、ブログの更新に応じて、ブログのWebページとして、1又は複数のブログ記事が掲載されるブログページを生成又は更新する。そして、ブログサーバ1は、ブログページDB101を備

10

20

30

40

50

え、ブログページを当該ブログページDB101に登録する。

【0031】

また、ブログサーバ1は、システム管理者により指定されたブロガーのブログページに対して、広告コンテンツ(関連コンテンツの一例)を挿入する。広告コンテンツとしては、例えば、広告文章のテキストデータ、バナーの画像データ、動画データ、Adobe Flash(商標)やSilverlight(商標)等により生成されたリッチインターネットアプリケーション(RIA)等がある。各ブログページに挿入される広告コンテンツは、対象のブログページに掲載されているブログ記事に関連する商品やサービスに関する広告を示すコンテンツである。そのため、ブログサーバ1は、複数の広告コンテンツが登録されている広告DB102を備える。そして、ブログサーバ1は、ブログページからブログ記事を抽出し、更にブログ記事から特徴語を抽出して、抽出した特徴語に関連する広告コンテンツを選択する。

10

【0032】

ユーザ端末3は、ブロガーとしてのユーザや、ブログを閲覧するユーザにより利用される端末装置である。ユーザ端末3としては、例えば、パーソナルコンピュータ、PDA、携帯電話機等が用いられる。

【0033】

管理端末2は、ブログシステムSのシステム管理者により使用される端末装置である。管理端末2としては、例えば、パーソナルコンピュータ等が用いられる。

【0034】

[2. ブログサーバの構成及び機能]

次に、ブログサーバ1の構成及び機能について、図2を用いて説明する。

20

【0035】

図2は、本実施形態に係るブログサーバ1の概要構成の一例を示すブロック図である。また、図3は、ブロガーが指定されてからブログページに広告コンテンツが挿入されるまでの処理の概要を示す図である。また、図4は、Webページの構成例を示す図である。また、図5は、HTML文書から生成されたDOMツリーの一例を示す図である。また、図6は、記憶部15に記憶されたコンテンツブロック対応情報の内容の一例を示す図である。

【0036】

図2に示すように、ブログサーバ1は、操作部11と、表示部12と、通信部13と、ドライブ部14と、記憶手段の一例としての記憶部15と、入出力インタフェース部16と、システム制御部20と、を備えている。そして、システム制御部20と入出力インタフェース部16とは、システムバス21を介して接続されている。

30

【0037】

操作部11は、例えば、キーボード、マウス等により構成されており、システム管理者等からの操作指示を受け付け、その指示内容を指示信号としてシステム制御部20に出力するようになっている。表示部12は、例えば、CRT(Cathode Ray Tube)ディスプレイ、液晶ディスプレイ等により構成されており、文字や画像等の情報を表示するようになっている。通信部13は、ネットワークNW等に接続して、管理端末2、ユーザ端末3等との通信状態を制御するようになっている。ドライブ部14は、例えば、フレキシブルディスク、CD(Compact Disc)、DVD(Digital Versatile Disc)等のディスクDKからデータ等を読み出す一方、当該ディスクDKに対してデータ等を記録するようになっている。

40

【0038】

記憶部15は、例えば、ハードディスクドライブ等により構成されており、各種プログラム及びデータ等を記憶するようになっている。また、記憶部15には、ブログページDB101及び広告DB102が構築されている。ブログページDB101には、ブログサービスサイトを構成する各ブログページ(ブログページのHTML文書(ドキュメントデータの一例)、ブログページの素材である画像データ等)が、例えば、そのページのUR

50

L及びブロガーの識別情報であるユーザIDに対応付けて登録されている。広告DB102には、複数の広告コンテンツが、その広告コンテンツによる広告対象の商品やサービスに関連するキーワードに対応付けて登録されている。ここで、広告コンテンツにテキストデータ以外のコンテンツが含まれている場合には、当該コンテンツのURLも対応付けて登録されている。また、ブログページ上に表示された広告コンテンツをユーザが選択したときに広告対象の商品やサービスに関するWebページに移動するため、当該WebページのURLも対応付けて登録されている。

【0039】

入出力インタフェース部16は、操作部11～記憶部15とシステム制御部20との間のインタフェース処理を行うようになっている。システム制御部20は、CPU(Central Processing Unit)17、ROM(Read Only Memory)18、RAM(Random Access Memory)19等により構成されている。

10

【0040】

システム制御部20は、CPU17が、ROM18や記憶部15に記憶された各種プログラムを読み出し実行することによりブログサーバ1の各部を制御する。また、システム制御部20は、広告コンテンツ挿入ソフトウェア(特有コンテンツ判定プログラムの一例)を実行することにより、抽出手段、計算手段、判定手段及び挿入手段として機能する。なお、広告コンテンツ挿入ソフトウェア等は、例えば、他のサーバ装置等からネットワークNWを介して取得されるようにしても良いし、CD-ROM等のディスクDKに記録されてドライブ部14を介して読み込まれるようにしても良い。

20

【0041】

広告コンテンツ挿入ソフトウェアは、ブログページに広告コンテンツを挿入するためのプログラムである。図3に示すように、広告コンテンツ挿入ソフトウェアは、マネージャ部、素材抽出エンジン、文章解析エンジン、広告選択部等により構成されている。マネージャ部は、素材抽出エンジン、文章解析エンジン及び広告選択部の実行を制御する。素材抽出エンジンは、ブログページのHTML文書からWeb素材としてのコンテンツを抽出するとともに、ブログページ特有のコンテンツを判定するためのソフトウェアである。コンテンツの抽出は、後述のコンテンツブロック(コンテンツグループの一例)の単位で行われる。本実施形態においては、例えば、その記事特有の内容を含むブログ記事が、ブログページに特有のコンテンツブロックに相当する。

30

【0042】

文章解析エンジンは、ブログページ特有のコンテンツとして抽出されたブログ記事から当該ブログページの特徴語を抽出するためのソフトウェアである。広告選択部は、抽出された特徴語をキーワードとして、ブログページに関連する広告コンテンツを選択するためのソフトウェアである。

【0043】

以下に、広告コンテンツの挿入の概要について説明する。図3に示すように、システム管理者により対象のブロガーのユーザIDが指定される(1)。すると、システム制御部20は、指定されたユーザIDに対応する全ブログページのHTML文書をブログページDB101から取得して解析し、Web素材たるコンテンツをコンテンツブロック単位で抽出する。そして、その抽出結果として、抽出したコンテンツブロック毎にコンテンツブロック対応情報(コンテンツ情報の一例)を生成する(2)。次いで、システム制御部20は、抽出した各コンテンツブロックの、指定されたユーザIDに対応する全ブログページにおける出現頻度を計算する。本実施形態において計算される出現頻度は、例えば、出現回数(度数)である。そして、システム制御部20は、各ブログページにおいて、出現頻度が所定の閾値以下であるコンテンツブロックを、そのブログページ特有のコンテンツブロックであると判定する(3)。

40

【0044】

システム制御部20は、特有のコンテンツブロックであると判定したコンテンツブロック、すなわち、ブログ記事に対して形態素解析等の分析を行って、ブロックページ毎の特

50

徴語を抽出する(4)。特徴語の抽出方法としては、種々の方法があり、また公知であるので、詳細な説明は省略する。一例としては、出現頻度の最も高い単語を特徴語とする。

【0045】

次いで、システム制御部20は、広告DB102を参照し、抽出した特徴語に関連する広告コンテンツを選択する(5)。そして、システム制御部20は、選択した広告コンテンツをブログページに挿入して表示させるための規定(タグやデータそのものの記述等)を、ブログページのHTML文書に挿入する(6)。

【0046】

次に、コンテンツブロックの抽出方法について説明する。なお、本実施形態においては、テキストデータ及び画像データがWeb素材として抽出されるものとする。

10

【0047】

例えば、ブログページの構成(レイアウト)が図4に示すようなものであるとする。Web素材としての各コンテンツは、ブログページ上において、それぞれ或るまとまり(かたまり)毎に表示されている。その各まとまりがコンテンツブロックに相当する。各コンテンツは、HTML文書に記述されているDIVタグ及びTABLEタグ(予め定められたタグの一例)により夫々コンテンツブロックに分けられる。つまり、各コンテンツは、DIVタグ及びTABLEタグ(以下、「ブロック化タグ」という)によりブロック化(グループ化)される。

【0048】

図4には、コンテンツブロック301~306が表示されている。コンテンツブロック301は、例えば、ページのヘッダ部分のコンテンツブロックであり、テキストA及び画像aにより構成されている。また、コンテンツブロック302は、例えば、各種カテゴリーの商品に関するWebページに移動するためのナビゲーション部分のコンテンツブロックであり、例えば、他のWebページへのリンクを示すテキストB、テキストC及びテキストDにより構成されている。また、コンテンツブロック303は、例えば、ブログの表示領域に対応するコンテンツブロックであり、ブログ等の見出しを示すテキストE、コンテンツブロック304及びコンテンツブロック305により構成されている。このように、コンテンツブロックが入れ子状、つまり、階層構造になっている場合もある。この場合、コンテンツブロック303に含まれるコンテンツはテキストEのみとされ、コンテンツブロック304及びコンテンツブロック305は、コンテンツブロック303から独立しているものとされる。コンテンツブロック304及び305は、夫々1件のブログ記事である。コンテンツブロック304には、ブログ記事の表題や本文を示すテキストF及びGにより構成されている。コンテンツブロック305には、ブログ記事の表題や本文を示すテキストH、I及びJと、ブログ記事に関連してブロガーにより登録された画像b及びcとにより構成されている。コンテンツブロック306は、例えば、コピーライト表示を示すコンテンツブロックであり、テキストIにより構成されている。

20

30

【0049】

これらのコンテンツブロックのうち、コンテンツブロック301、302、303及び306は、図4に示すブログページ以外のブログページ上でも比較的頻繁に現れる。一方、コンテンツブロック304及びコンテンツブロック305は、基本的に当該ブログページにだけに用いられる。よって、コンテンツブロック304又はコンテンツブロック305が、当該ブログページ特有のコンテンツブロックであると判断されることとなる。

40

【0050】

本実施形態においては、その記事特有の内容を含むブログ記事に相当するコンテンツブロックを特有のコンテンツブロックと判定されるようにする必要がある。こうした特有の内容を含むブログ記事は、1ページ内に複数含まれている場合がある。そのため、出現頻度が所定の閾値以下であるコンテンツブロックは、全て特有のコンテンツブロックとされる。例えば、閾値の値を1回に設定する。そうすると、特有の内容を含むブログ記事は、特有のコンテンツブロックと判定され、他のブログ記事と同じような内容のみを含むブログ記事は、特有のコンテンツブロックとは判定されない。また、ヘッダ部分、ナビゲーシ

50

ョン部分、コピーライト表示部分等の各ブログページに共通するようなコンテンツブロックの出現頻度は夫々2回以上となるので、これらも特有のコンテンツブロックとは判定されない。なお、閾値は記憶部15に予め記憶される。

【0051】

図4に示すブログページのHTML文書をDOM(Document Object Model)ツリー、すなわち、木構造で表したものが図5である。なお、図5に示すDOMツリーにおいて、本実施形態の説明に必要なタグのノードの図示は省略している。

【0052】

DOMツリーにおいては、DIVタグを示すDIVノード及びTABLEタグを示すTABLEノードが、各コンテンツをコンテンツブロックにブロック化するノードとされる(以下、「ブロック化ノード」という)。システム制御部20は、例えば、深さ優先探索によりDOMツリーを探索し、コンテンツブロックを確定していく。具体的に、システム制御部20は、ブロック化ノードを発見すると、当該ノードを頂点とする部分木の各ノードに規定されているコンテンツを一まとめにしてコンテンツブロックとする。ただし、コンテンツブロックが階層的に規定された結果、或るブロック化ノード(以下、「上位のブロック化ノード」という)を発見した後、当該ブロック化ノードの子孫のノードの中から更にブロック化ノード(以下、「下位のブロック化ノード」という)を発見すると、コンテンツブロックが分割される。例えば、ルートノードからの距離が近いノードほど階層的により上位のノードであるとする、上位のブロック化ノードを頂点とする部分木(以下、「上位部分木」という)に相当するコンテンツブロックを、下位のブロック化ノードを頂点とする部分木(以下、「下位部分木」という)に相当するコンテンツブロックと、上位部分木のうち下位部分木を除く部分に相当するコンテンツブロックとに分ける(例えば、コンテンツブロック304とコンテンツブロック303)。この場合においては、前者のコンテンツブロックを、後者のコンテンツブロックよりも階層的に下位にあるものとする。例えば、コンテンツブロック301、302、303及び306の階層は1であり、コンテンツブロック304及び305の階層は2である。つまり、階層の値が小さいほど階層的により上位となる。

【0053】

これを、HTML文書のタグの記述で説明すると、基本的にはブロック化タグが記述されていると、ブロック化タグに挟まれた範囲内に規定が記述されているコンテンツが、まとめて当該ブロック化タグに対応するコンテンツブロックとされる。ただし、ブロック化タグが階層的に記述されている場合、或るブロック化タグに挟まれた範囲内に規定が記述されているコンテンツのうち、当該ブロック化タグよりも下位に記述されているブロック化タグに挟まれた範囲に規定が記述されているコンテンツを除いたものが、まとめて当該ブロック化タグに対応するコンテンツブロックとされる。

【0054】

このようにしてコンテンツブロックを抽出すると、システム制御部20は、抽出結果を示すコンテンツブロック対応情報を一時的に記憶部15に保存する。図6に示すように、コンテンツブロック対応情報(符号401)は、コンテンツブロック毎に保存される。コンテンツブロック対応情報には、抽出元のHTML文書のURL設定部分(符号402)とブロック構成情報(符号403)とにより構成されている。ブロック構成情報には、抽出された各コンテンツが設定される。ここで、テキストデータについては、DOMツリーにおけるテキストノードの内容が設定される。一方、画像データについては、画像データそのものの代わりとして、DOMツリーにおいてIMGタグを示すIMGノードのsrc属性としての画像データのURLが設定される。なお、本実施形態においては、ブログページ特有のコンテンツブロックであると判定されたコンテンツブロック、すなわち、ブログ記事から特徴語を抽出するので、テキストデータが抽出されれば良く、画像データについては抽出しなくても良い。

【0055】

[3. ブログシステムの動作]

10

20

30

40

50

次に、ブログシステム S の動作について、図 7 乃至図 10 を用いて説明する。

【0056】

図 7 は、本実施形態に係るブログサーバ 1 のシステム制御部 20 の広告コンテンツ挿入処理における処理例を示すフローチャートである。

【0057】

広告コンテンツ挿入処理は、例えば、システム管理者の操作に基づいて管理端末 2 から広告コンテンツ挿入処理の実行のリクエストが送信されてきたときに開始される。

【0058】

そして、システム管理者が広告コンテンツの挿入対象とするブログを運営するブロガーのユーザ ID を指定すると、図 7 に示すように、システム制御部 20 は、指定されたユーザ ID を管理端末 2 から受信する（ステップ S 1）。

10

【0059】

次いで、システム制御部 20 は、ブロック数 NUM に 0 を設定する（ステップ S 2）。ブロック数 NUM は、現時点で発見済みのコンテンツブロックの個数である。また、NUM はグローバル変数であり、後述する 1 ページ対応抽出処理及びツリー探索処理からアクセスが可能である。

【0060】

次いで、システム制御部 20 は、受信したユーザ ID に対応する最初のブログページの HTML 文書をブログページ DB 101 から取得する（ステップ S 3）。次いで、システム制御部 20 は、取得した HTML 文書を指定して、後述する 1 ページ対応抽出処理を実行する（ステップ S 4）。この 1 ページ対応抽出処理では、取得した HTML 文書からコンテンツブロックが抽出され、コンテンツブロック対応情報が保存される。

20

【0061】

次いで、システム制御部 20 は、受信したユーザ ID に対応する全てのブログページのコンテンツブロックを抽出したか否かを判定する（ステップ S 5）。このとき、システム制御部 20 は、コンテンツブロックを抽出していないブログページが存在する場合には（ステップ S 5：NO）、次のブログページの HTML 文書をブログページ DB 101 から取得して（ステップ S 6）、ステップ S 4 に移行する。そして、システム制御部 20 は、ステップ S 4～S 6 の処理を繰り返して全てのブログページのコンテンツブロックを抽出すると（ステップ S 5：YES）、ステップ S 7 に移行する。

30

【0062】

ステップ S 7 において、システム制御部 20 は、受信したユーザ ID に対応する最初のブログページの HTML 文書を特定する。

【0063】

次いで、システム制御部 20 は、取得した HTML 文書を指定して、後述する特有コンテンツブロック判定処理を実行する（ステップ S 8）。この特有コンテンツブロック判定処理では、特定した HTML 文書からコンテンツブロックが抽出され、ブログページ特有のコンテンツブロックが判定される。

【0064】

次いで、システム制御部 20 は、特有と判定されたコンテンツブロックを構成する各テキストデータからブログページの特徴語を抽出する（ステップ S 9）。次いで、システム制御部 20 は、抽出した特徴語に基づいて、ブログページに関連する広告ページを当該ブログページに挿入する（ステップ S 10）。具体的に、システム制御部 20 は、抽出した特徴語をキーワードとし、広告 DB 102 に参照して当該キーワードに対応する広告コンテンツを選択する。次いで、システム制御部 20 は、特定した HTML 文書上の所定の位置に、選択した広告コンテンツの規定を挿入する。例えば、システム制御部 20 は、広告コンテンツにテキストデータが含まれている場合には、当該テキストデータの内容を HTML 文書に追加する。また、例えば、システム制御部 20 は、広告コンテンツに画像データが含まれている場合には、当該画像データを表示するための IMG タグを HTML 文書に追加する。また、例えば、システム制御部 20 は、広告対象の商品やサービスに関する

40

50

Web ページへのリンク情報を HTML 文書に追加する。

【 0 0 6 5 】

システム制御部 2 0 は、特定した HTML 文書に広告コンテンツの規定を挿入すると、当該 HTML 文書で、ブログページ DB 1 0 1 に登録されている HTML 文書を更新する (ステップ S 1 1)。

【 0 0 6 6 】

次いで、システム制御部 2 0 は、受信したユーザ ID に対応する全てのブログページに広告コンテンツを挿入したか否かを判定する (ステップ S 1 2)。このとき、システム制御部 2 0 は、広告コンテンツを挿入していないブログページが存在する場合には (ステップ S 1 2 : NO)、次のブログページの HTML 文書を特定して (ステップ S 1 3)、ステップ S 8 に移行する。そして、システム制御部 2 0 は、ステップ S 8 ~ S 1 3 の処理を繰り返して全てのブログページに広告コンテンツを挿入すると (ステップ S 1 2 : YES)、記憶部 1 5 に保存させておいた全てのコンテンツブロック対応情報を、記憶部 1 5 から削除する (ステップ S 1 4)。システム制御部 2 0 は、この処理を終えると、広告コンテンツ挿入処理を終了させる。

10

【 0 0 6 7 】

図 8 は、本実施形態に係るブログサーバ 1 のシステム制御部 2 0 の 1 ページ対応抽出処理における処理例を示すフローチャートである。

【 0 0 6 8 】

図 8 に示すように、システム制御部 2 0 は、先ず、取得した HTML 文書の DOM ツリーを RAM 1 9 上に生成する (ステップ S 2 1)。

20

【 0 0 6 9 】

次いで、システム制御部 2 0 は、階層 LV に 0 を設定する (ステップ S 2 2)。階層 LV は、DOM ツリーにおいて現在探索中のノードが属するコンテンツブロックの階層である。LV はグローバル変数であり、1 ページ対応抽出処理及び後述するツリー探索処理からアクセスが可能である。

【 0 0 7 0 】

次いで、システム制御部 2 0 は、DOM ツリーのルートノードを指定して (ステップ S 2 3)、ツリー探索処理を実行する (ステップ S 2 4)。ツリー探索処理は再帰呼び出しが可能であり、このツリー探索処理により、Web ページから全てのコンテンツブロックが抽出され、コンテンツブロック対応情報が生成される。

30

【 0 0 7 1 】

次いで、システム制御部 2 0 は、ツリー探索処理により生成された各コンテンツブロック対応情報を記憶部 1 5 に保存する (ステップ S 2 5)。システム制御部 2 0 は、この処理を終えると、1 ページ対応抽出処理を終了させる。

【 0 0 7 2 】

図 9 は、本実施形態に係るブログサーバ 1 のシステム制御部 2 0 のツリー探索処理における処理例を示すフローチャートである。

【 0 0 7 3 】

図 9 に示すように、システム制御部 2 0 は、先ず、指定されたノードの種類を判定する (ステップ S 3 1)。このとき、システム制御部 2 0 は、指定されたノードの種類が、DIV ノード又は TABLE ノード (ブロック化ノード) である場合には、すなわち、コンテンツブロックが発見された場合には (ステップ S 3 1 : DIV 又は TABLE)、ステップ S 3 2 に移行する。

40

【 0 0 7 4 】

ステップ S 3 2 において、システム制御部 2 0 は、ブロック数 NUM に 1 を加算するとともに、階層 LV に 1 を加算する。次いで、システム制御部 2 0 は、ブロック番号 BN [LV] に NUM を設定する (ステップ S 3 3)。ブロック番号 BN [LV] は、現在探索中のノードが属する階層 LV で示されるコンテンツブロックのブロック番号である。このブロック番号は、コンテンツブロックの発見順に付与される。また、BN [LV] は、グ

50

ローバル変数である。

【0075】

次いで、システム制御部20は、ブロック番号BN[L V]のコンテンツブロックに対応するコンテンツブロック対応情報を初期化する(ステップS34)。具体的に、システム制御部20は、コンテンツブロック対応情報を格納する領域をRAM19上に設定し、取得したHTML文書のURLを、当該領域に設定する。

【0076】

次いで、システム制御部20は、指定されたノードの子ノードのうち、未だ探索されていない子ノードが存在するか否かを判定する(ステップS35)。このとき、システム制御部20は、未だ探索されていない子ノードが存在する場合には(ステップS35:YES) 10、ステップS36に移行する。

【0077】

ステップS36において、システム制御部20は、探索されていない子ノードのうちの1つの子ノードを指定して、ツリー探索処理を実行する(ステップS37)。システム制御部20は、ツリー探索処理を終えると、ステップS35に移行する。

【0078】

そして、システム制御部20は、ステップS35~S37の処理を繰り返して全ての子ノードのツリー探索処理を終えると(ステップS35:NO)、ステップS38に移行する。なお、システム制御部20は、指定されたノードの子ノードが1つも存在しない場合にも、ステップS38に移行する。ステップS38において、システム制御部20は、階層LVから1を減算して、ツリー探索処理を終了させる。 20

【0079】

ステップS31において、システム制御部20は、指定されたノードの種類がテキストノードである場合には(ステップS31:テキスト)、指定されたノードの内容(テキストデータ)を、ブロック番号BN[L V]のコンテンツブロックに対応するコンテンツブロック対応情報中のブロック構成情報に追加設定する(ステップS39)。システム制御部20は、この処理を終えると、ツリー探索処理を終了させる。

【0080】

ステップS31において、システム制御部20は、指定されたノードの種類がIMGノードである場合には(ステップS31:IMG)、指定されたノードのsrc属性として 30 設定されている画像データのURLを取得し、取得したURLを、ブロック番号BN[L V]のコンテンツブロックに対応するコンテンツブロック対応情報中のブロック構成情報に追加設定する(ステップS40)。システム制御部20は、この処理を終えると、ツリー探索処理を終了させる。

【0081】

ステップS31において、システム制御部20は、指定されたノードの種類が、DIVノード、TABLEノード、テキストノード、及びIMGノードの何れでもない場合には(ステップS31:その他)、指定されたノードの子ノードのうち、未だ探索されていない子ノードが存在するか否かを判定する(ステップS41)。このとき、システム制御部20は、未だ探索されていない子ノードが存在する場合には(ステップS41:YES) 40、探索されていない子ノードのうちの1つの子ノードを指定して(ステップS42)、ツリー探索処理を実行する(ステップS43)。システム制御部20は、ツリー探索処理を終えると、ステップS41に移行する。

【0082】

一方、システム制御部20は、指定されたノードの全ての子ノードのツリー探索処理を終えた場合、又は、指定されたノードの子ノードが1つも存在しない場合には(ステップS41:NO)、ツリー探索処理を終了させる。

【0083】

図10は、本実施形態に係るコンテンツ生成サーバ1のシステム制御部20の特有コンテンツブロック判定処理における処理例を示すフローチャートである。 50

【 0 0 8 4 】

図 1 0 に示すように、システム制御部 2 0 は、先ず、1 ページ対応抽出処理と同様に、指定された HTML 文書の DOM ツリー生成 (ステップ S 6 1)、ブロック数 NUM 及び階層 LV に対して 0 の設定を行い (ステップ S 6 2)、DOM ツリーのルートノードを指定して (ステップ S 6 3)、ツリー探索処理を実行する (ステップ S 6 4)。

【 0 0 8 5 】

次いで、システム制御部 2 0 は、ブロック番号 i に 1 を設定する (ステップ S 6 5)。次いで、システム制御部 2 0 は、ブロック番号 i のコンテンツブロックの出現頻度を計算する (ステップ S 6 6)。

【 0 0 8 6 】

具体的に、システム制御部 2 0 は、ステップ S 6 4 のツリー探索処理において生成されたコンテンツブロック対応情報 i (ブロック番号 i のコンテンツブロックに対応するコンテンツブロック対応情報) のブロック構成情報と、記憶部 1 5 に保存されている各コンテンツブロック対応情報のブロック構成情報とを比較する。このとき、システム制御部 2 0 は、ブロック構成情報の内容が一致する場合には、出現回数 1 回としてカウントする。このとき、システム制御部 2 0 は、ブロック構成情報中におけるコンテンツの規定順は無視してかまわない。また、システム制御部 2 0 は、記憶部 1 5 に保存されているコンテンツブロック対応情報のブロック構成情報に規定されている一部のコンテンツがコンテンツブロック対応情報 i のブロック構成情報に規定されている全部のコンテンツに一致する場合も、出現回数 1 回としてカウントしても良い。更に、システム制御部 2 0 は、コンテンツ
20
ブロック対応情報のブロック構成情報中に規定されているテキストデータ同士を比較する場合には、テキストデータが示す文章等そのものが一致するか否かを判定するのではなく、その文章等により表現されている実質的な内容を比較しても良い。例えば、システム制御部 2 0 は、夫々のテキストデータの形態素解析等を行うことによりテキストデータから単語を抽出し、抽出した単語同士を比較しても良い。そして、システム制御部 2 0 は、全ての単語が一致した場合にテキストデータ同士が一致したと判断しても良いし、所定の割合以上で単語が一致した場合にテキストデータ同士が一致したと判断しても良い。システム制御部 2 0 は、このようにしてコンテンツブロック対応情報 i のブロック構成情報と、記憶部 1 5 に保存されている全てのコンテンツブロック対応情報のブロック構成情報とを比較して、出現頻度を計算する。
30

【 0 0 8 7 】

システム制御部 2 0 は、出現頻度を計算すると、計算した出現頻度が、記憶部 1 5 に記憶されている閾値以下であるか否かを判定する (ステップ S 6 7)。このとき、システム制御部 2 0 は、出現頻度が閾値以下である場合には (ステップ S 6 7 : YES)、ブロック番号 i のコンテンツブロックを、特有のコンテンツブロックの 1 つであると判定する (ステップ S 6 8)。つまり、システム制御部 2 0 は、ブロック番号 i のコンテンツブロックを、指定された HTML 文書が対応するブログページに特有のコンテンツブロックに加える。

【 0 0 8 8 】

システム制御部 2 0 は、出現頻度が閾値よりも大きい場合 (ステップ S 6 7 : NO)、
40
又は、ステップ S 6 8 の処理を終えた場合には、ブロック番号 i に 1 を加算して (ステップ S 6 9)、ブロック番号 i がブロック数 NUM の値より大きいか否かを判定する (ステップ S 7 0)。このとき、システム制御部 2 0 は、ブロック番号 i がブロック数 NUM の値以下である場合には (ステップ S 7 0 : NO)、ステップ S 6 6 に移行する。そして、システム制御部 2 0 は、ツリー探索処理において抽出された全てのコンテンツブロックの出現頻度を計算すると (ステップ S 7 0 : YES)、特有コンテンツブロック判定処理を終了させる。

【 0 0 8 9 】

なお、システム制御部 2 0 は、ステップ S 6 4 のツリー探索処理によりコンテンツブロックを抽出していたが、広告コンテンツ挿入処理から実行された 1 ページ対応抽出処理 (50

図7ステップS4)において、受信したブロガーのユーザIDに対応する全てのブログページについてコンテンツブロックが抽出され、その結果としてコンテンツブロック対応情報が記憶部15に記憶されているので、再度コンテンツブロックを抽出しなくても良い。その場合には、指定されたHTML文書のURLに基づいて、当該HTML文書が対応するブログページを構成する各コンテンツブロックのコンテンツブロック対応情報を記憶部15から取得することができる。

【0090】

[4.変形例1]

次に、本実施形態の変形例について、図11を用いて説明する。

【0091】

これまでの説明においては、システム管理者によりブロガーが指定されたときに、指定されたブロガーのブログページに広告コンテンツを挿入していたが、ブログが更新されたタイミングで広告コンテンツを挿入しても良い。

【0092】

図11は、本実施形態の変形例に係るブログサーバ1のシステム制御部20のブログ更新時処理における処理例を示すフローチャートである。なお、図11において、図7と同様の処理については同様のステップ番号を付してある。

【0093】

まず、ブログの更新に先立ち、ブロガーは、ユーザ端末3を操作してブログサービスサイトにアクセスし、自身のユーザIDとパスワードとを入力することによりブログサービスサイトにログインする。このログインにより、ブログサーバ1はユーザ端末3に対してセッションIDを発行し、セッションIDとユーザIDとを対応付けて管理する。ユーザ端末3からブログサーバ1へのリクエストにはセッションIDが含まれているので、ブログサーバ1は、どのブロガーからのリクエストであるかを特定することができる。

【0094】

そして、ブロガーが新しいブログ記事の登録操作を行うと、ユーザ端末3は、ブログ記事のデータ(表題や本文等のテキストデータ、画像データ等)をブログサーバ1に送信し、図11に示すように、ブログサーバ1のシステム制御部20は、ブログ記事のデータを受信する(ステップS71)。次いで、システム制御部20は、ブロガーのユーザIDに対応するブログページの中から、更新すべきブログページのHTML文書をブログページDB101から取得する(ステップS72)。次いで、システム制御部20は、受信したブログ記事のデータに基づいて、取得したHTML文書を更新する(ステップS73)。例えば、システム制御部20は、取得したHTML文書に、ブログ記事用のTABLEタグ又はDIVタグを追加し、当該タグに挟まれた形で、受信したブログ記事の表題や本文のテキストデータ等を追加する。次いで、システム制御部20は、ブログ記事のデータを追加したHTML文書で、ブログページDB101に登録されているHTML文書を更新する(ステップS74)。

【0095】

次いで、システム制御部20は、ブロガーのユーザIDに対応する全てのブログページからコンテンツブロックを抽出する(ステップS3~S6)。

【0096】

次いで、システム制御部20は、ステップS73において更新したHTML文書を指定して、特有コンテンツブロック判定処理を実行し(ステップS8)、特有と判定されたコンテンツブロックを構成する各テキストデータからブログページの特徴語を抽出する(ステップS9)。

【0097】

次いで、システム制御部20は、指定されたHTML文書から、既存の広告コンテンツの規定を削除し(ステップS75)、抽出した特徴語をキーワードとして、関連する広告コンテンツの規定を挿入する(ステップS10)。つまり、システム制御部20は、ブログページ上に表示される広告コンテンツを変更する。

10

20

30

40

50

【 0 0 9 8 】

そして、システム制御部 2 0 は、広告コンテンツの規定が挿入された HTML 文書で、ブログページ DB 1 0 1 に登録されている HTML 文書を更新し（ステップ S 1 1）、記憶部 1 5 から全てのコンテンツブロック対応情報を削除する（ステップ S 1 4）。

【 0 0 9 9 】

なお、ブログの更新に伴ってブログページを新規に生成しなければならない場合の処理も、基本的に上述した処理と同様で良い。ただし、新規に生成されたブログページには、広告コンテンツは未だ挿入されていないので、ステップ S 7 5 における広告コンテンツの規定の削除は行われぬ。

【 0 1 0 0 】

[5 . 変形例 2]

これまでの説明においては、ブログページに特有のコンテンツの判定に用いられる閾値として 1 回を設定していたが、2 回以上の値を閾値として設定しても良い。

【 0 1 0 1 】

例えば、閾値を 1 回とした場合には、出現頻度が 1 回であるコンテンツブロック（ブログ記事）がブログページに特有のコンテンツとして抽出され、抽出されたブログ記事のテキストデータから特徴語が抽出される。このとき、抽出された各ブログ記事のテキストデータのデータ量が少ないと、そこから抽出される単語の数は少なくなる。そして、十分な数の単語を抽出することができず、どの単語が特徴語であるかを全く判断することができない場合や、的確に判断することができない場合がある。そこで、閾値の値を上げて、ブログページに特有のコンテンツと判定される条件をゆるめることで、特徴語を抽出する対象となるブログ記事を増やしていく。これにより、特徴語を抽出することが可能となる。

【 0 1 0 2 】

具体的には、ブログサーバ 1 のシステム制御部 2 0 が、最初は閾値を 1 回に設定して、ブログページに特有のコンテンツブロックを判定することにより、出現回数が 1 回のブログ記事を抽出して特徴語を抽出する。このとき、システム制御部 2 0 は、特徴語を抽出することができずと判定した場合には、閾値を 2 回に変更して、ブログ記事の抽出及び特徴語の抽出を行う。システム制御部 2 0 は、それでも特徴語を抽出することができずと判定した場合には、閾値を 3 回に変更して、ブログ記事の抽出及び特徴語の抽出を行う。システム制御部 2 0 は、こうした処理を、特徴語が抽出することができるまで継続する。つまり、特有のコンテンツブロックの抽出結果に基づく処理を正常に行うことができなかった場合に、閾値を上げるのである。

【 0 1 0 3 】

ただし、閾値を無制限に上げていくと、ブログ記事ではないものも抽出されてしまうので、閾値がある程度まで上がると処理を中断するものとする。例えば、閾値が、指定されたブロガーに対応するブログページのページ数の値にまで上がると、各ブログページで共通して用いられるコンテンツブロックを抽出してしまうので、閾値がブログページのページ数の値になったら処理を中断しても良い。

【 0 1 0 4 】

また例えば、ブログページの所定ページ数あたり 1 回のみ出現するコンテンツブロックをブログページ特有のコンテンツブロックであると、システム管理者側で予め定めても良い。この場合、指定されたブロガーに対応するブログページのページ数に比例して、閾値としての出現回数を変えても良い。

【 0 1 0 5 】

[6 . 変形例 3]

これまでの説明においては、ブログページに特有のコンテンツの判定に用いられる出現頻度として、出現回数（度数）を用いていたが、指定されたブロガーに対応するブログページの全コンテンツブロックに対する出現回数の割合（相対度数）を用いても良い。

【 0 1 0 6 】

例えば、ブロガーが登録したブログ記事に対して、他のユーザからコメントを登録することができ、ブログ記事とともにコメントが閲覧可能になっているとする。このコメントのテキストデータもブログページを構成するコンテンツの1つとなる。ブログサーバ1のシステム制御部20は、コメントのテキストデータをブログページに追加する場合、当該ブログページのHTML文書に、ブロック化タグの記述を追加した上で当該テキストデータを追加することにより、コメントのテキストデータを、ブログ記事や他のコメントのテキストデータとは独立したコンテンツブロックとする。そして、システム制御部20は、コンテンツブロックとしてコメントのテキストデータを抽出し、抽出したコメントのテキストデータが特有の内容を有している場合には、そのコメントに関連する広告コンテンツをブログページに挿入するようにする。

10

【0107】

ところで、或るブログ記事に対して複数のコメントが登録された場合において、複数のコメントの各内容が、例えば、多数派の意見と少数派の意見といったように、頻繁に出現する内容と、あまり頻繁には出現しない内容とに分かれる場合がある。このとき、多数派の意見は、一般的な意見であり、あまり特徴的な内容ではないと考えることができる。一方、少数派の意見は、特異な意見であり、ブログページに特有の内容と考えることができる。そうした場合に、少数派の意見を示すコメントをブログページ特有のコンテンツとして抽出したい。

【0108】

しかしながら、多数派の意見の数と少数派の意見の数は、相対的なものであり、コメントの総数によって変化する。こうした場合において、出現頻度として度数を用い、閾値を例えば1回とすると、頻繁に現れない内容(少数派の意見)を適切に抽出することができない場合がある。そこで、出現頻度として相対度数を用い、閾値を所定の割合に設定するのである。このときの閾値は任意に設定することができる。例えば、抽出されたコンテンツブロックの内容がN個のパターン(Nは2以上の整数)に分けられる場合、少数派の意見を区別するために、閾値には $1 \div N$ 未満の範囲で閾値を設定しても良い。このように、システム制御部20が、その時々状況に応じて閾値を変更しても良い。

20

【0109】

なお、ブログ等のような記事に対してコメント等を登録することができるシステムとして、例えば、或るユーザが登録したつぶやきに対して、これをフォローするつぶやきを他のユーザが登録することができるTwitter(商標)や、電子掲示板等がある。

30

【0110】

以上説明したように、本実施形態によれば、ブログサーバ1のシステム制御部20が、HTML文書が指定されることによって順次指定されたブログページを構成しているコンテンツを抽出し、指定されたブログページを構成している各コンテンツの出現頻度を計算し、指定されたブログページを構成するコンテンツのうち、出現頻度が所定の閾値以下のコンテンツを当該ブログページに特有のコンテンツであると判断する。

【0111】

従って、出現頻度が小さいコンテンツであるほど、指定されたブログページ以外にはあまり出現しないコンテンツであるので、出現頻度が閾値以下であるかを判定することで、当該条件を満たす全てのコンテンツが、指定されたブログページに特有のコンテンツであると特定される。よって、ブログページに特有のコンテンツを容易に抽出することができる。

40

【0112】

また、ブログサーバ1のシステム制御部20が、指定されたブログページに特有のコンテンツに関連する広告コンテンツを当該ブログページに挿入する。

【0113】

従って、ブログページの特徴と関連する情報をWebページに追加することができる。

【0114】

また、ブログサーバ1のシステム制御部20が、指定されたブログページを構成してい

50

るコンテンツとして、ブログの記事のテキストデータが含まれている場合に、当該テキストデータを、当該ブログページに特有のコンテンツであると判定し、ブログの記事のテキストデータから当該ブログページの特徴語を抽出し、当該特徴語をキーワードとして、予め関連付けられている広告コンテンツを、当該ブログページに挿入する。

【0115】

従って、ブログページに掲載されているブログの内容に関連する広告を当該ブログページに追加することができる。

【0116】

また、ブログサーバ1のシステム制御部20が、ブログサービスサイトに含まれる複数のブログページ上における各コンテンツの出現頻度を計算する。

10

【0117】

従って、ブログサービスサイトに含まれる複数のWebページ(例えば、指定されたブロガーのユーザIDに対応する複数のブログページ)上において、指定されたブログページを構成している各コンテンツの出現頻度が計算されるので、ブログサービスサイト内で共通して用いられるコンテンツは、特有のコンテンツではないと判定することが可能となり、判断精度を上げることができる。

【0118】

また、ブログサーバ1のシステム制御部20が、1つ以上のコンテンツで構成されるコンテンツブロックの単位で、ブログページを構成しているコンテンツを抽出し、指定されたブログページを構成している各コンテンツブロック出現頻度を計算し、指定されたブログページを構成するコンテンツブロックのうち、出現頻度が閾値以下のコンテンツブロックを当該ブログページに特有のコンテンツブロックであると判断する。

20

【0119】

従って、ブログページ上において、例えば、ヘッダ部分、ナビゲーション部分、ブログが表示される部分、コピーライト表示の部分等のように、1つ以上のコンテンツがまとまりをもってコンテンツブロックとして表示されている場合に、ブログページに特有のコンテンツブロックを抽出することができる。

【0120】

また、ブログサーバ1のシステム制御部20が、ブログページを構成しているコンテンツを当該ブログページのHTML文書に基づいて抽出し、HTML文書においてDIVタグ又はTABLEタグに基づいて、コンテンツブロックを定める。

30

【0121】

従って、DIVタグにより、HTML文書の作成の際に明示的にブロック化された1つ以上のコンテンツを特定することができ、また、TABLEタグにより、表形式でブロック化されて表示される1つ以上のコンテンツを特定することができるので、例えば、これらのタグにより、ブログページに特有のコンテンツと、特有ではないコンテンツとがブロック化されている場合に、Webページに特有のコンテンツを判断する精度を上げることができる。

【0122】

なお、上記実施形態においては、Webページを構成しているコンテンツとして、テキストデータ及び画像データを抽出していたが、抽出対象のコンテンツはこれらに限られるものではない。例えば、Webページ上に表示されるコンテンツ、又は、Webページが表示されている際に再生されるコンテンツ(例えば、動画データ、音声データ、電子文書等)であれば良い。また、所定の種類のコンテンツのみを抽出しても良い。

40

【0123】

また、上記実施形態においては、指定されたブログページを構成する各コンテンツブロックに対応するコンテンツブロック対応情報を、指定されたブロガーのユーザIDに対応する全てのブログページを構成する各コンテンツブロック対応情報と比較することによって各出現頻度が計算されていた。つまり、指定されたブログページを構成する各コンテンツブロックの出現頻度を計算する場合に、指定されたブロガーに対応する全てのブログペ

50

ージを対象とした範囲に出現する頻度を計算するようになっていたが、対象とする範囲はこれだけに限られるものではない。例えば、予め定められたページ数分のブログページを対象としても良いし、ブログサービスサイトを構成する全てのブログページを対象としても良い。

【0124】

また、上記実施形態においては、D I Vタグに挟まれているコンテンツ、及び、T A B L Eタグに挟まれているコンテンツを、コンテンツブロックとしてグループ化して抽出していたが、コンテンツをグループ化するタグとしては、これらのみに限られるものではない。

【0125】

また、上記実施形態においては、W e b ページに特有のコンテンツをコンテンツブロックの単位で抽出していたが、各コンテンツをそのまま一つずつ抽出しても良い。

【0126】

また、上記実施形態として、W e b ページに特有のコンテンツに関連するコンテンツとして、商品やサービスに関する広告を示す広告コンテンツを、当該W e b ページに挿入していたが、関連するコンテンツであれば広告コンテンツに限られるものではない。例えば、特有のコンテンツと判定されたブログ記事等のコンテンツに関連する画像データ（静止画や動画像）を、背景画像や挿入画像（挿絵等）として挿入しても良い。具体的には、例えば、画像データ用のデータベースを構築し、当該データベースに、画像データとキーワードとを対応付けて登録しておく。画像データに対応付けられるキーワードは、その画像データによって表される画像を示す単語や当該画像に関連する単語である。そして、特有のコンテンツと判定されたコンテンツから特徴語を抽出し、抽出した特徴語をキーワードとして、関連する画像データをデータベースから選択する。そして、対象のH T M L 文書のB O D Yタグに、選択した画像データのU R Lをb a c k g r o u n d属性として挿入したり、対象のH T M L 文書の所定位置に、選択した画像データを表示するI M Gタグを挿入したりする。これにより、特有のコンテンツと判定されたブログ記事等のコンテンツの内容に適した画像をW e b ページに挿入することができる。

【0127】

また、W e b ページに特有のコンテンツの用途としては、関連するコンテンツをW e b ページに挿入することのみに限られるものではない。例えば、W e b ページに特有のコンテンツに基づいて、新たなコンテンツを生成しても良い。

【0128】

また、上記実施形態においては、サーバ装置に対して本発明の特有コンテンツ判定装置を適用していたが、記憶手段やネットワーク上からH T M L 文書を取得することができれば、端末装置等に対して特有コンテンツ判定装置を適用しても良い。

【0129】

また、上記実施形態においては、H T M L 文書に対して本発明のドキュメントデータを適用していたが、マークアップ言語で記述され、W e b ページを構成するコンテンツを示すデータ（例えば、X H T M L（Extensible HyperText Markup Language）文書等）に対してドキュメントデータを適用しても良い。

【0130】

また、上記実施形態においては、ブログサービスサイトにおけるブログページを構成するコンテンツを抽出していたが、対象とするサイト及びページの種類はこれらのみに限られるものではない。

【符号の説明】

【0131】

- 1 ブログサーバ
- 2 管理端末
- 3 ユーザ端末
- 11 操作部

10

20

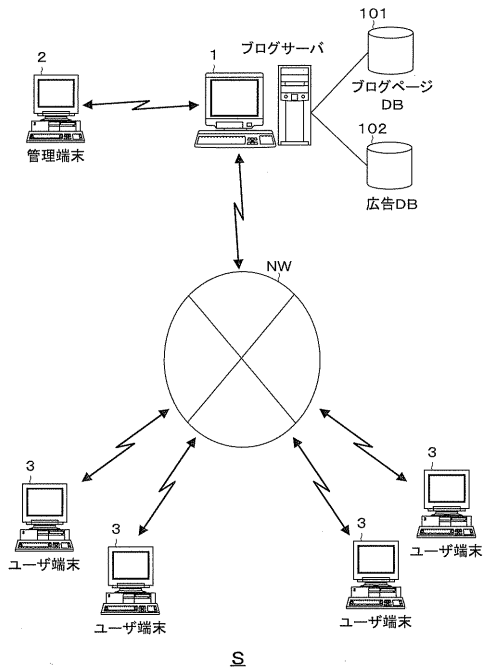
30

40

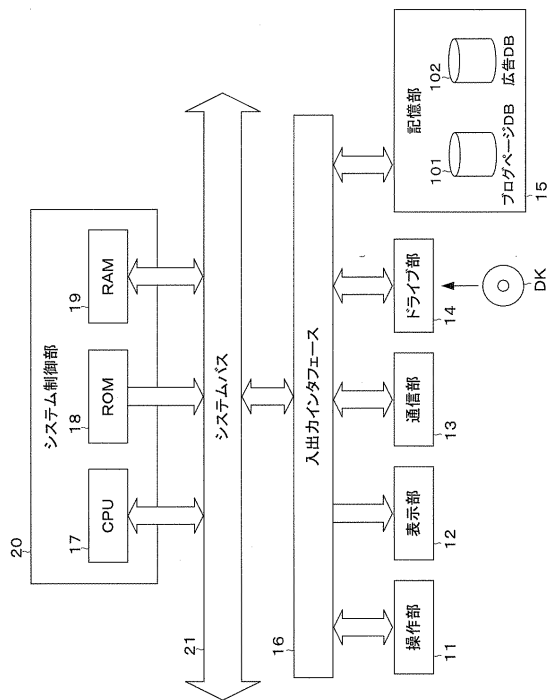
50

- 1 2 表示部
- 1 3 通信部
- 1 4 ドライブ部
- 1 5 記憶部
- 1 6 入出インタフェース部
- 1 7 CPU
- 1 8 ROM
- 1 9 RAM
- 2 0 システム制御部
- 2 1 システムバス
- 1 0 1 ブログページDB
- 1 0 2 広告DB
- NW ネットワーク
- S ブログシステム

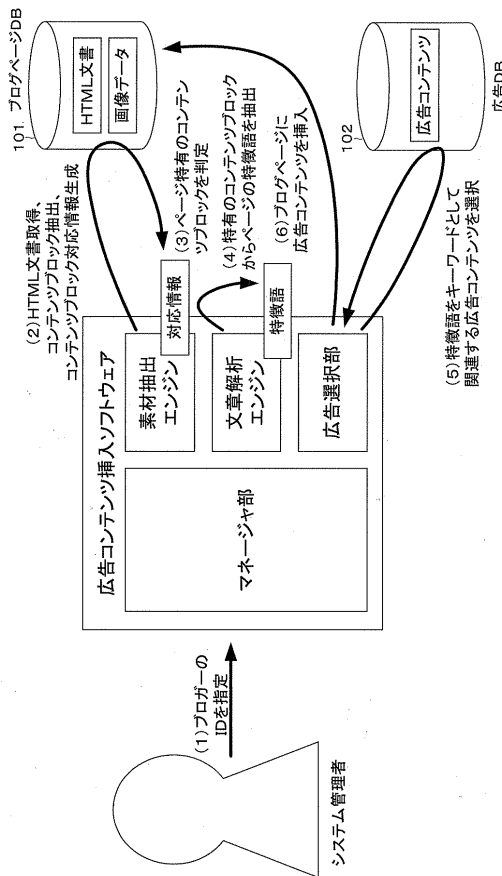
【図1】



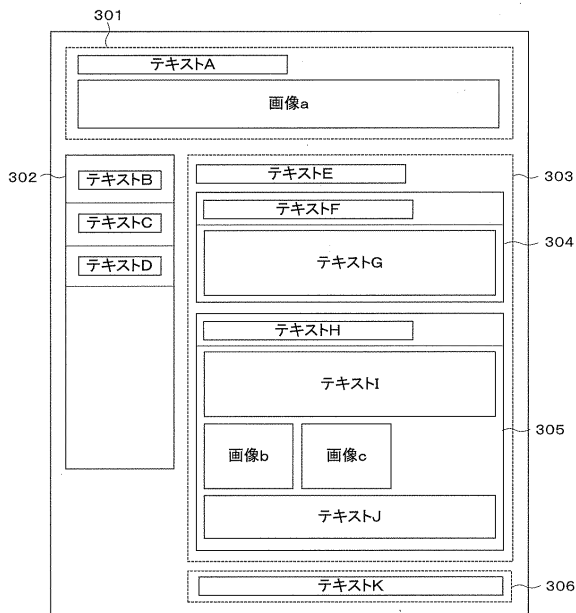
【図2】



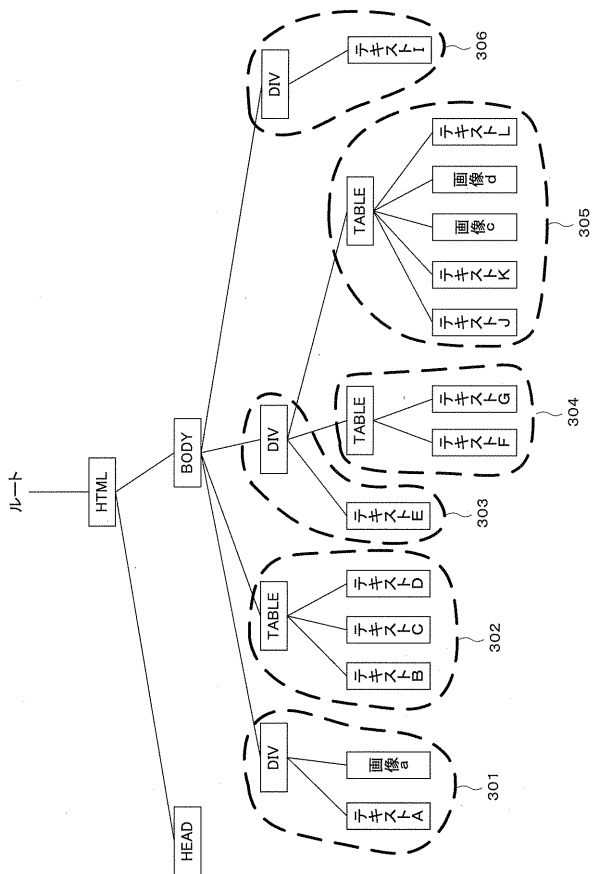
【図3】



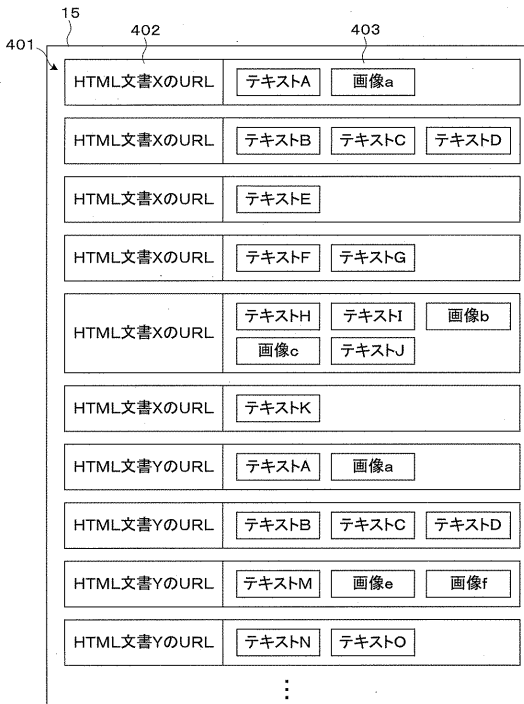
【図4】



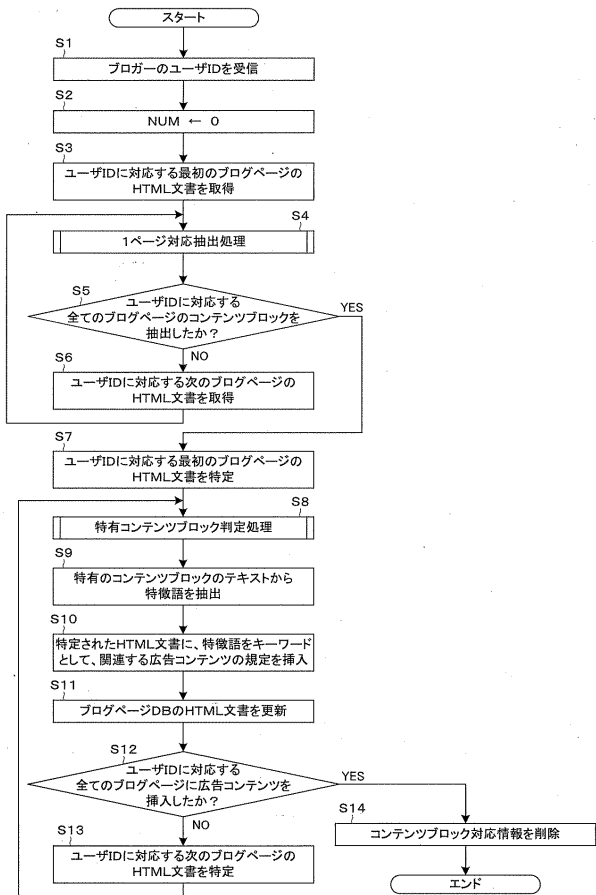
【図5】



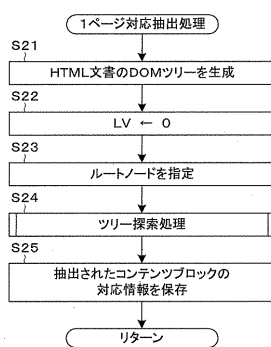
【図6】



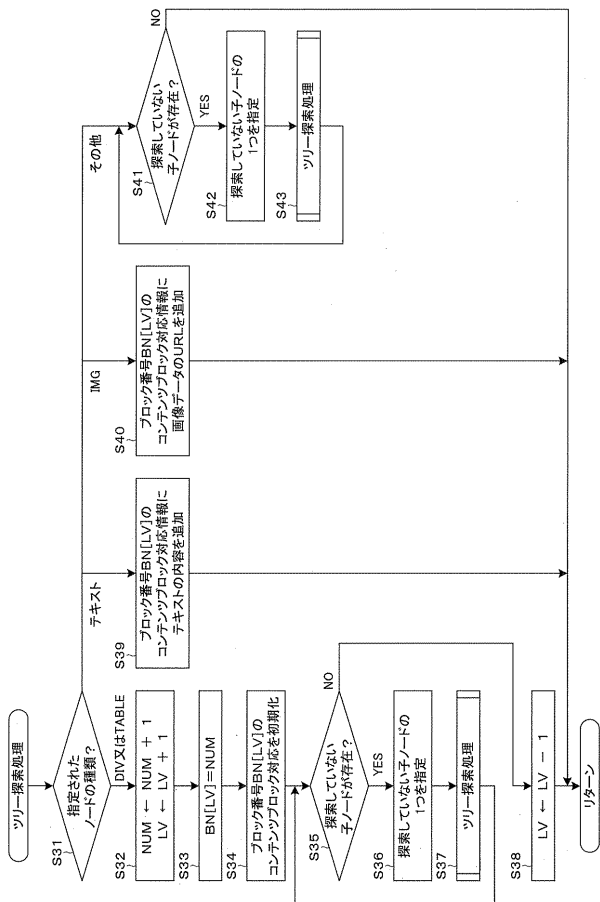
【図7】



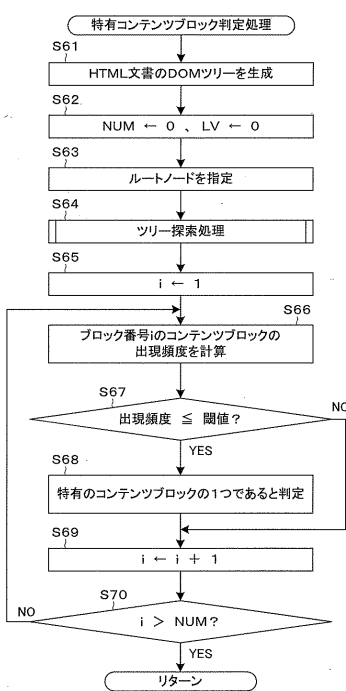
【図8】



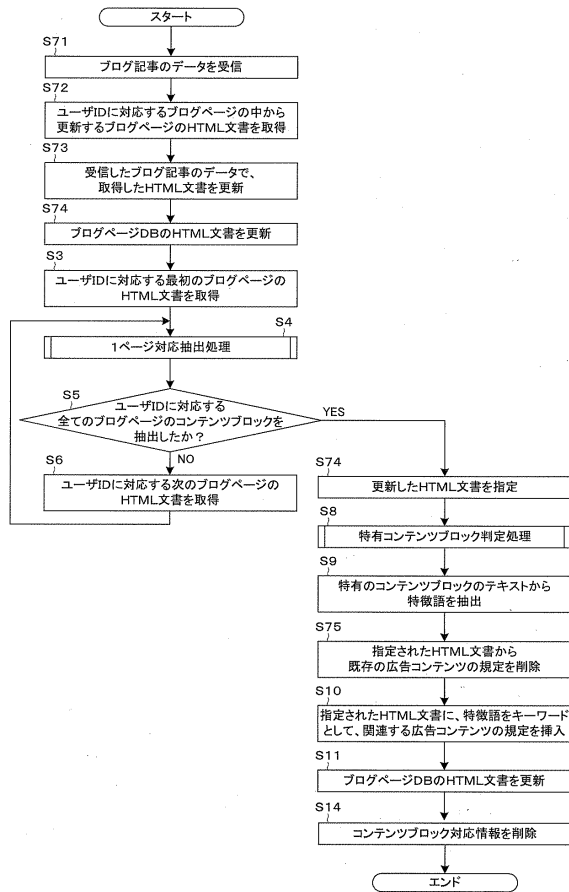
【図9】



【図10】



【図11】



フロントページの続き

- (56)参考文献 特開2005-302041(JP,A)
国際公開第2008/108515(WO,A1)
特表2010-533897(JP,A)
特開2005-122690(JP,A)
特開2008-139928(JP,A)
特開2009-053983(JP,A)
特開2009-205499(JP,A)
特開2006-146506(JP,A)
Sandip Debnath, Prasenjit Mitra, Nirmal Pal, and C. Lee Giles, "Automatic Identification of Informative Sections of Web Pages", IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 17, NO. 9, SEPTEMBER 2005, pp. 1233-1246
吉田光男, 山本幹雄, "教師情報を必要としないWebページ群のコンテンツ自動抽出ツールの提案", Internet<URL:http://db-event.jpn.org/deim2009/proceedings/files/A8-4.pdf>, 第1回データ工学と情報マネジメントに関するフォーラム-DEIMフォーラム-論文集, DEIMフォーラム2009, 電子情報通信学会データ工学研究専門委員会, 2009年05月09日

- (58)調査した分野(Int.Cl., DB名)
G06Q10/00-50/00