

(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2012-138027

(P2012-138027A)

(43) 公開日 平成24年7月19日(2012.7.19)

(51) Int.Cl. F I テーマコード(参考)
G06F 17/30 (2006.01) G06F 17/30 340Z 5B075
 G06F 17/30 170A

審査請求 有 請求項の数 9 O L (全 20 頁)

(21) 出願番号 特願2010-291331(P2010-291331)
 (22) 出願日 平成22年12月27日(2010.12.27)

(71) 出願人 000003078
 株式会社東芝
 東京都港区芝浦一丁目1番1号
 (71) 出願人 301063496
 東芝ソリューション株式会社
 東京都港区芝浦一丁目1番1号
 (74) 代理人 100149803
 弁理士 藤原 康高
 (72) 発明者 藤田 慎一
 東京都港区芝浦一丁目1番1号 東芝ソ
 リューション株式会社内
 (72) 発明者 高知尾 勝彦
 東京都港区芝浦一丁目1番1号 東芝ソ
 リューション株式会社内

最終頁に続く

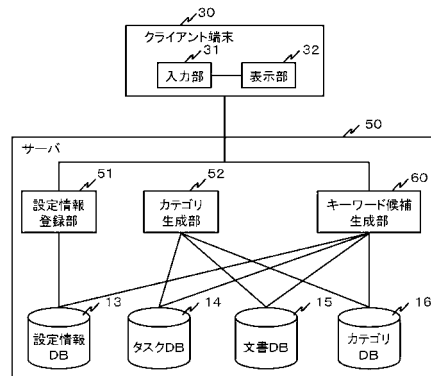
(54) 【発明の名称】 情報検索システム、検索キーワード提示方法、およびプログラム

(57) 【要約】 (修正有)

【課題】 単語間の関連を保持する辞書や、ユーザが入力した検索キーワードと参照文書を記録するログ情報などを用いず、所望する情報を検索可能な検索キーワードの候補をユーザに提示できる情報検索システムを提供する。

【解決手段】 情報検索システムは、検索対象の文書データを格納する文書データベースと検索対象の文書データに対するクラスタリング分析によるカテゴリ分類に関する情報を含むカテゴリ情報を格納するカテゴリデータベースとを備える。また、情報検索システムは、ユーザによって入力される検索キーワードを含む検索条件に基づいて文書データベースから文書データを検索し、検索した文書データがカテゴリのうちどのカテゴリに属するかを判定し、検索した文書データが属すると判定されたカテゴリ、もしくは当該カテゴリの文書データから、キーワード候補を抽出し、抽出したキーワード候補をユーザに提示する。

【選択図】 図2



【特許請求の範囲】**【請求項 1】**

検索対象の文書データを格納する文書データベースと、
カテゴリに前記検索対象の文書データのうちのどの文書データが属するかを管理する
カテゴリ情報を格納するカテゴリデータベースと、
ユーザによって入力される検索キーワードを含む検索条件に基づいて前記文書データベ
ースから文書データを検索する文書検索部と、
前記カテゴリデータベースを参照して、前記文書検索部によって検索された文書デー
タが前記カテゴリのうちのどのカテゴリに属するかを判定するカテゴリ判定部と、
前記カテゴリ判定部によって前記検索された文書データが属すると判定されたカテゴリ
、もしくは、当該カテゴリに属する文書データから、キーワード候補を抽出するキーワ
ード候補抽出部と、
前記キーワード候補抽出部によって抽出された前記キーワード候補をユーザに提示する
出力部と、
を備える情報検索システム。

10

【請求項 2】

前記キーワード候補抽出部は、前記カテゴリ判定部によって前記検索された文書デー
タが属すると判定されたカテゴリの文書データに含まれる単語のうち、前記検索キーワ
ード以外の単語の出現頻度を計算し、前記出現頻度があらかじめ設定された所定数以上である
前記単語を前記キーワード候補として抽出する請求項 1 に記載の情報検索システム。

20

【請求項 3】

前記検索条件は前記文書データが前記文書データベースに登録された日時の範囲を指定
するための検索期間情報を含み、
前記キーワード候補抽出部は、前記カテゴリ判定部によって前記検索された文書デー
タが属すると判定されたカテゴリの文書データから、前記検索期間情報に基づいて文書デー
タの抽出を行い、当該抽出された文書データにおいて、検索キーワード以外の単語におけ
る出現頻度を計算し、前記出現頻度があらかじめ設定された所定数以上である前記単語を
前記キーワード候補として抽出する請求項 1 乃至請求項 2 のいずれか 1 項に記載の情報検
索システム。

30

【請求項 4】

前記キーワード候補抽出部は、前記カテゴリ判定部によって前記検索された文書デー
タが属すると判定されたカテゴリの名称を前記キーワード候補として抽出する請求項 1 乃至
請求項 3 のいずれか 1 項に記載の情報検索システム。

【請求項 5】

前記キーワード候補抽出部は、前記カテゴリ判定部によって前記検索された文書デー
タが属すると判定されたカテゴリに含まれる文書データに対して前記クラスタリング分析を
行い、

当該クラスタリング分析の結果、新たに作成されたカテゴリが、前記検索された文書デー
タを含む場合に、当該作成されたカテゴリの名称を前記キーワード候補として抽出する
請求項 1 乃至請求項 4 のいずれか 1 項に記載の情報検索システム。

40

【請求項 6】

前記キーワード候補抽出部は、前記検索された文書データが属する前記カテゴリに含ま
れる文書データから、前記検索された文書データと関連付けられた文書データを抽出し、
当該抽出された文書データから、キーワード候補を抽出する請求項 1 乃至請求項 4 のい
ずれか 1 項に記載の情報検索システム。

【請求項 7】

前記文書検索部は、あらかじめ定められた周期が到来した場合に、前記文書データを検
索する請求項 1 乃至請求項 6 のいずれか 1 項に記載の情報検索システム。

【請求項 8】

ユーザによって入力される検索キーワードを含む検索条件に基づいて文書データベ
ース

50

から文書データを検索するステップと、

カテゴリデータベースを参照して、検索された前記文書データがどのカテゴリに属するかを判定するステップと、

検索された前記文書データが属すると判定されたカテゴリ、もしくは、当該カテゴリに属する文書データから、キーワード候補を抽出するステップと、

抽出された前記キーワード候補をユーザに提示するステップと、

を備える検索キーワード提示方法。

【請求項 9】

検索対象の文書データを格納する文書データベースと、前記検索対象の文書データに対するクラスタリング分析によるカテゴリ分類に関する情報を含むカテゴリ情報を格納するカテゴリデータベースとを備える情報検索システムのプログラムであって、

コンピュータに、

ユーザによって入力される検索キーワードを含む検索条件に基づいて文書データベースから文書データを検索する機能と

前記検索された文書データが前記カテゴリのうちのどのカテゴリに属するかを判定する機能と、

前記検索された文書データが属すると判定されたカテゴリ、もしくは、当該カテゴリの文書データから、キーワード候補を抽出する機能と、

抽出した前記キーワード候補をユーザに提示する機能と、を実現させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は情報検索システム、検索キーワード提示方法、およびプログラムに関する。

【背景技術】

【0002】

従来、文書が大量に蓄積されているデータベースよりユーザが所望している情報を見つける検索装置がある。

【0003】

この検索装置において、検索に使用するキーワードはユーザが指定するものであるため、キーワードを指定するユーザ自身が知らない単語は検索キーワードとして使用することができない。これにより、ユーザは必要とする情報を取得できない場合がある。もしくは、ユーザが検索で使用するキーワードが廃れた場合にも、ユーザは必要とする情報を取得できない場合がある。

【0004】

これらの問題に対し、単語と単語に関連する関連単語とを対応付ける関連辞書などを用いて検索を実現する装置がある。また、複数のカテゴリに分類した文書情報に対して検索を行い、参照要求した文書の識別情報をログ情報から取得することで、ユーザが入力した検索キーワードとカテゴリの適合度から検索キーワードを提示することを可能とする装置もある。

【0005】

しかしながら、上述したような従来技術では、検索キーワードと関連付いている単語を探すような場合に、単語間の関連性に関する情報を保持する関連辞書が必要である。そのため、辞書を作成するための労力や費用の発生、新たに出てくる単語を辞書に適用するといったメンテナンスが発生するというような問題がある。

【0006】

また、カテゴリに分類された文書の中からユーザが検索で用いたキーワードを元にどの文書を参照したかを知るにはログ情報とログ情報の解析が必要である。したがって、このログ情報を取得する手段と、取得したログ情報を解析する手段が必要になるという問題がある。

10

20

30

40

50

【先行技術文献】

【特許文献】

【0007】

【特許文献1】特開2008-70991号公報

【発明の概要】

【発明が解決しようとする課題】

【0008】

本発明が解決しようとする課題は、単語間の関連を保持する辞書、又はユーザが入力した検索キーワードと参照した文書を記録するログ情報を用いずに、ユーザが所望する情報を検索可能な検索キーワードの候補をユーザに提示することを可能とする情報検索システム、検索キーワード提示方法、およびプログラムを提供することである。

10

【課題を解決するための手段】

【0009】

実施形態の情報検索システムは、検索対象の文書データを格納する文書データベースと、検索対象の文書データに対するクラスタリング分析によるカテゴリ分類に関する情報を含むカテゴリ情報を格納するカテゴリデータベースとを備える。実施形態の情報検索システムは、ユーザによって入力される検索キーワードを含む検索条件に基づいて文書データベースから文書データを検索し、検索された文書データがカテゴリのうちどのカテゴリに属するかを判定し、検索された文書データが属すると判定されたカテゴリ、もしくは、当該カテゴリの文書データから、キーワード候補を抽出し、抽出したキーワード候補をユーザに提示する。

20

【図面の簡単な説明】

【0010】

【図1】実施形態の情報検索システムの全体構成を示す図。

【図2】実施形態の情報検索システムの機能構成を示す図。

【図3】実施形態の情報検索システムの検索条件設定画面の一例を示す図。

【図4】実施形態の情報検索システムの処理手順の一例を示すフローチャート。

【図5】実施形態の情報検索システムのキーワード候補生成部の機能ブロック図。

【図6】実施形態の情報検索システムのキーワード候補生成部によるキーワード候補生成処理の一例を示すフローチャート。

30

【図7】実施形態の情報検索システムのカテゴリ解析結果の概念の一例を示す図。

【図8】実施形態の情報検索システムのカテゴリ名取得部によるキーワード候補抽出処理の一例を示すフローチャート。

【図9】実施形態の情報検索システムの新カテゴリ名取得部によるキーワード候補抽出処理の一例を示すフローチャート。

【図10】実施形態の情報検索システムの新カテゴリ名取得部によるキーワード候補抽出処理の概念図。

【図11】実施形態の情報検索システムの出現頻度解析部によるキーワード候補抽出処理の一例を示すフローチャート。

【図12】実施形態の情報検索システムの出現頻度解析部によるHTMLカテゴリからのキーワード候補抽出処理の概念図。

40

【図13】実施形態の情報検索システムの期間限定出現頻度解析部によるキーワード候補抽出処理の一例を示すフローチャート。

【図14】実施形態の情報検索システムのKL連携解析部によるキーワード候補抽出処理の一例を示すフローチャート。

【図15】実施形態の情報検索システムのKL連携解析部によるHTMLカテゴリからのキーワード候補抽出処理の概念図。

【図16】実施形態の情報検索システムの検索条件設定画面の一例を示す図。

【発明を実施するための形態】

【0011】

50

以下、本実施形態の情報検索システムについて図面を参照して説明する。

【0012】

図1は、本実施形態に係る検索キーワード提示方法を用いた情報検索システムの構成例を示すブロック図である。

【0013】

本実施形態の情報検索システムは、図1に示すように、クライアント端末からの入力に基づいて提示キーワードを作成するサーバコンピュータ(以下、「サーバ」という)50にローカルエリアネットワーク(Local Area Network: LAN)等のネットワーク20を介してクライアントコンピュータ(以下、「クライアント端末」という)30が複数台接続されたサーバクライアントシステムを想定する。サーバ50およびクライアント端末30は、例えば、一般的なパーソナルコンピュータである。

10

【0014】

サーバ50はコンピュータ10と記憶装置11を備える。コンピュータ10は、ハードディスクドライブなどの記憶装置11と接続されている。この記憶装置11は、コンピュータ10によって実行されるアプリケーションプログラムである本実施形態の情報検索システムのプログラム12を記憶している。すなわち、記憶装置11は、プログラム12を記憶する記憶媒体として機能する。

【0015】

クライアント端末30上では、コンピュータ10を利用するクライアントソフトウェアが動作する。クライアントソフトウェアは、例えばブラウザである。複数のクライアント端末30は、ネットワーク20を介してコンピュータ10と接続されている。なお、図1には、2台のクライアント端末30以外のクライアント端末30は省略されている。

20

【0016】

図2は本実施形態の情報検索システムの機能構成の一例を示すブロック図である。

【0017】

図2に示すように、本実施形態の情報検索システムのクライアント端末30は、入力部31と表示部32とを備える。入力部31は例えばキーボードやマウスであり、表示部32は例えば液晶ディスプレイである。

【0018】

ユーザは入力部31を用いて、ユーザが注目している文書(検索したい文書)に関するキーワード(以下、検索キーワードという)や、ユーザが文書データを検索する検索範囲などを含む検索条件をサーバに入力する。

30

【0019】

なお、ユーザによる設定情報の入力は、例えば、表示部32に表示される、文書データを検索するための検索条件を設定する画面である検索条件設定画面を用いて行う。

【0020】

入力部31から入力された検索条件はサーバ50が有する設定情報データベース13に設定情報として登録される。サーバ50は、この設定情報に基づいて周期的に文書検索を行い、検索結果の文書に基づいて、ユーザに検索キーワードの候補(以下、キーワード候補という)を提示する。具体的には、表示部32であるディスプレイに、この検索結果の文書に関するキーワード候補が表示される。すなわち、表示部32はキーワード候補をユーザに提示するための出力部として機能する。なお、サーバ50が文書検索もしくはキーワード候補抽出を行う周期は、あらかじめ設定される。

40

【0021】

本実施形態の情報検索システムのサーバ50は、設定情報登録部51と、カテゴリ生成部52と、キーワード候補生成部60とを備えて構成される。本実施形態において、設定情報登録部51、カテゴリ生成部52、およびキーワード候補生成部60はそれぞれ、図1に示されるコンピュータ10が記憶装置11に格納されているプログラム12を実行することにより実現されるものとする。このプログラム12は、コンピュータ読み取り可能な記憶媒体に予め格納して頒布可能である。また、このプログラム12が、ネットワーク

50

20を介してコンピュータ10にダウンロードされても構わない。

【0022】

また、サーバ50は、設定情報データベース(以下、DBという)13、タスクDB14、文書DB15、及びカテゴリDB16を含む。

【0023】

設定情報DB13には、ユーザが入力した検索条件が設定情報として記憶されている。設定情報は例えば、設定情報の識別子、検索キーワードなどを含む。

【0024】

ここで、図3を参照して、設定情報の入力、及び設定情報の登録について説明する。図3は検索条件設定画面100の一例である。

【0025】

図3に示すように、検索条件設定画面100は、「設定条件名」101、「検索キーワード」102、「推奨キーワード」103、「キーワード検索範囲」104、「タスク名」105、および「期間」106の項目を有する。図3の検索条件設定画面100では、設定条件名101は「コンピュータ最新動向」、検索キーワード102は、「インターネット」もしくは「アーキテクチャ」、キーワード検索範囲104は、「タイトル」、タスク名105は「最新ネットワーク技術動向」、期間106は、「3日前から」と入力されている。

【0026】

ユーザがこれらの項目を入力し、「設定ボタン」107を入力部31であるマウスでクリックすると、サーバ50にこれらの入力内容が入力される。入力された情報は、後述する設定情報登録部51によって設定情報として設定情報DB13に登録される。なお、本実施形態では、設定情報の識別子をユーザが入力する設定条件名101としているが、設定情報登録部51によって、入力された設定情報の識別子である設定情報IDが発行されるとしてもよい。

【0027】

タスクDB14には、タスクの識別子であるタスクID、タスク名、タスクに関する説明などを含むタスク情報が記憶されている。なお、タスクとは、クラスタリング分析(以下、クラスタリングという)によりカテゴリを生成する単位である。すなわち、タスクDB14に格納されているタスク情報はクラスタリング分析を行う際に用いられる。なお、クラスタリングを実行する際の対象となる文書データは、ある1つのタスクに属しているため、あるカテゴリに属している文書は全て同じタスクに属している文書となる。

【0028】

文書DB15には、文書の識別子である文書ID、文書のタイトル、文書データの本文、文書のナレッジリンク(以下、KLという)情報、文書の投稿日時、文書の属しているタスクのIDなどを含む文書データが記憶されている。

【0029】

なお、KL情報とは、文書DB15に格納されている複数の記事の相互間の関係を示す。すなわち、KLとは異なる2つの文書間の関係を表す情報である。なお、KLは、1つのタスク内に属している文書間のみでなく、別々のタスク内に属している文書間の情報としても生成される。すなわち、あるタスク以外のタスクに属している文書に含まれるキーワードを取得するには、あるタスクに属している文書とKLにより結びついている他のタスクに属している文書を取得し、その文書からキーワードを取得する。なお、文書DB15には、そのほかに、文書の作成者名を含んでもよいし、あらかじめユーザがラベルを付与しても良い。

【0030】

すなわちKL情報は、例えばKL情報毎に文書データの相互間の関係の元となる文書を識別する元文書ID、当該文書相互間の関係の先となる文書を識別する先文書ID、及び当該元文書IDと先文書IDとの関係が示されている。

【0031】

10

20

30

40

50

カテゴリDB16は、例えばカテゴリ名などのカテゴリに関する情報と、当該カテゴリに属する文書データなどを含むカテゴリ情報が記憶されている。本実施形態の情報検索システムにおいては、カテゴリDB16に格納されたカテゴリは、後述するカテゴリ生成部52によって、タスクDB14と文書DB15を参照して作成されるとする。

【0032】

なお、本実施形態において、これらのDBは、記憶装置11に格納される。

【0033】

設定情報登録部51は、ユーザが検索条件設定画面にて設定した、文書を検索するための条件情報を設定情報として設定情報DB13に登録する。

【0034】

カテゴリ生成部52は、タスクDB14と文書DB15とを参照して、どのタスクにどの文書が属しているかという情報を取得する。

【0035】

そしてカテゴリ生成部52は、取得した情報に基づいて、タスク毎にクラスタリング分析を実施し、カテゴリを生成する。カテゴリ生成部52は、生成されたカテゴリに関する情報と、カテゴリにどの文書が属しているかという情報とをカテゴリ情報としてカテゴリDB16に登録する。なお、このカテゴリ生成処理はあらかじめ設定された期間ごとに定期的に行うものとする。

【0036】

キーワード候補生成部60は、設定情報DB13から取得した設定情報に基づいて、タスクDB14と文書DB15とから文書データを取得する。そして、キーワード候補生成部60は、カテゴリDB16のカテゴリ情報に基づいて、取得した文書データがどのカテゴリに属しているかを判定する。

【0037】

キーワード候補生成部60は、上述したカテゴリ判定の結果、取得された文書が特定数以上属しているカテゴリのカテゴリ名、取得された文書が特定数以上属し、かつカテゴリが新規に作成されたカテゴリのカテゴリ名、取得された文書が特定数以上属しているカテゴリ内の文書内において出現頻度の高い単語、取得された文書が特定数以上属しているカテゴリ内の文書であって、期間106によって指定された期間内に生成されている文書内において出現頻度の高い単語、又は、取得された文書が特定数以上属しているカテゴリ内で取得された文書とKLで関連付いている文書内において出現頻度の高い単語などをキーワード候補として抽出する。抽出されたキーワード候補は、表示部32の検索条件設定画面における推奨キーワード103欄に表示される。なお、抽出されるキーワード候補はユーザが入力した検索キーワード102以外の語句である。

【0038】

すなわち、キーワード候補生成部60は、ユーザが入力した検索キーワードに基づいて検索された文書データを検索可能なキーワードをキーワード候補として抽出する。これらの抽出されたキーワード候補は、ユーザが所望する情報を取得する精度を向上できる。

【0039】

ここで、図4を参照して、本実施形態の情報検索システムの動作について説明する。図4は、本実施形態の情報検索システムにおけるユーザに検索キーワードを提示する処理（以下、検索キーワード提示処理という）の手順の一例を示すフローチャートである。

【0040】

まず、入力部31によって、設定情報が入力される（ステップS1）。入力された設定情報は、ネットワーク20を介してクライアント端末30からサーバ50に送信される。

【0041】

サーバ50が設定情報を受信すると、設定情報登録部51が設定情報DB13に受信した設定情報を登録する（ステップS2）。

【0042】

あらかじめ設定された周期が到来すると（ステップS3がYes）、キーワード候補生

10

20

30

40

50

成部 60 が、設定情報 DB 13 に登録された設定情報に基づいて文書 DB 15 からキーワードを抽出し、キーワード候補を生成する（ステップ S4）。抽出されたキーワード候補の一覧は、表示部 32 に表示された検索条件設定画面 100 の推奨キーワード 103 に表示される（ステップ S5）。これにより、検索キーワード提示処理が終了する。なお、推奨キーワード 103 に表示されたキーワード候補をユーザが選択することで、選択されたキーワード候補が検索キーワード 102 に追加されるようにしても良い。また、この場合、追加されたキーワードは設定情報 DB 13 に更新されるようにしても良い。

【0043】

なお、あらかじめ設定された周期が到来していない場合は（ステップ S3 が No）、その周期が来るまで待機する。

【0044】

ここで、図 5 乃至図 15 を参照して、図 4 のステップ S4 におけるキーワード候補生成部 60 によるキーワード候補の生成（以下、キーワード候補生成処理という）について、詳しく説明する。

【0045】

図 5 はキーワード候補生成部 60 の機能ブロック図である。図 5 に示すように、キーワード候補生成部 60 は、注目文書検索部 61、カテゴリ判定部 62、キーワード候補抽出部 63 を備える。

【0046】

注目文書検索部 61 は、特定の周期が到来すると、設定情報 DB 13 より設定情報を取得し、タスク DB 14 と文書 DB 15 を参照して、設定情報 DB 13 に登録されている検索キーワード 102、及び検索対象の期間 106、検索対象のタスク名に合致する文書を抽出する。注目文書検索部 61 は、この抽出結果である文書データをカテゴリ判定部 62 に送信する。

【0047】

カテゴリ判定部 62 は、受信した文書のデータに基づいて、カテゴリ DB 16 よりカテゴリ情報を取得する。そして、カテゴリ判定部 62 は、取得したカテゴリ情報と、注目文書検索部 61 から受信した検索結果である文書の一覧とに基づいて、検索結果に含まれている当該文書それぞれがどのカテゴリに属しているかを解析する。解析結果である各文書がどのカテゴリに属しているかという情報はキーワード候補抽出部 63 に送信される。

【0048】

キーワード候補抽出部 63 は、出現頻度解析部 64、期間限定出現頻度解析部 65、カテゴリ名取得部 66、新カテゴリ名取得部 67、および KL 連携取得部 68 を備える。キーワード候補抽出部 63 は、カテゴリ判定部 62 から受信した解析結果に基づいて、キーワード候補を抽出する。そして、抽出されたキーワード候補の一覧は表示部 32 の検索条件設定画面 100 の推奨キーワード 103 に表示される。以下、キーワード候補抽出部 63 が備える、出現頻度解析部 64、期間限定出現頻度解析部 65、カテゴリ名取得部 66、新カテゴリ名取得部 67、および KL 連携取得部 68 について説明する。

【0049】

出現頻度解析部 64 は、キーワード抽出部 63 がカテゴリ判定部 62 から受信した解析結果に基づいて、検索結果の文書を含んでいる各カテゴリ内に属している文書における、検索キーワード 102 に設定されている単語以外で出現頻度の高い単語をキーワード候補として抽出する。出現頻度解析部 64 は、抽出したキーワード候補をキーワード候補抽出部 63 に送信する。

【0050】

期間限定出現頻度解析部 65 は、キーワード抽出部 63 がカテゴリ判定部 62 から受信した解析結果から、検索結果の文書を含んでいる各カテゴリ内に属していて、期間 106 によって指定された期間内で生成された文書において、検索キーワード 102 に設定されている単語以外で出現頻度の高い単語をキーワード候補として抽出する。期間限定出現頻度解析部 65 は、抽出したキーワード候補をキーワード候補抽出部 63 に送信する。

10

20

30

40

50

【 0 0 5 1 】

カテゴリ名取得部 6 6 は、キーワード抽出部 6 3 がカテゴリ判定部 6 2 から受信した解析結果から、カテゴリ内に含まれる文書数をカテゴリごとに計測する。カテゴリ名取得部 6 6 は計測結果に基づき、特定数以上の注目文書を含むカテゴリのカテゴリ名をキーワード候補として抽出する。カテゴリ名取得部 6 6 は、抽出したキーワード候補をキーワード候補抽出部 6 3 に送信する。

【 0 0 5 2 】

新カテゴリ名取得部 6 7 は、キーワード抽出部 6 3 がカテゴリ判定部 6 2 から受信した解析結果から、カテゴリ内の検索結果の文書数をカテゴリごとに計測する。新カテゴリ名取得部 6 7 は計測結果に基づき、特定数以上の検索結果の文書を含むカテゴリの中でクラスタリング分析を行う。そして、新カテゴリ名取得部 6 7 はクラスタリング分析により新たに作成されたカテゴリのなかで、検索結果の文書が含まれるカテゴリのカテゴリ名を検索キーワード候補として抽出する。新カテゴリ名取得部 6 7 は、抽出したキーワード候補をキーワード候補抽出部 6 3 に送信する。

10

【 0 0 5 3 】

K L 連携取得部 6 8 は、キーワード抽出部 6 3 がカテゴリ判定部 6 2 から受信した解析結果から、検索結果の文書を含んでいる各カテゴリ内に属していて、検索結果の文書と K L により関連付いている文書内において、検索キーワードに設定されている単語以外で出現頻度の高い単語を検索キーワード候補として抽出する。K L 連携取得部 6 8 は、抽出したキーワード候補をキーワード候補抽出部 6 3 に送信する。

20

【 0 0 5 4 】

なお、単語の頻出頻度の高さや、カテゴリが含む検索結果の文書の数の多さやを判定するための閾値は、あらかじめシステム構築時などに設定されているとする。また、ユーザが設定可能なようにしてもよい。また、閾値は文書 DB 1 5 に登録されたデータ量に対する検索結果の文書のデータ量の割合などでも良い。もしくは閾値を設定するのではなく、一番検索結果の文書が多く属するカテゴリを抽出するとしても良い。

【 0 0 5 5 】

また、K L とは、異なる 2 つの文書間の因果関係を表す情報であり、ある文書から見た他の文書への参照とその説明から成る。また、K L には因果関係に沿った方向性があり、入力としての元文書から出力としての先文書へという関係を保持している。K L は、例えば「Describe」、「Follow」、「Revise」、「Sum」、「Consult」という種類（以下、クラスという）を有する。

30

【 0 0 5 6 】

なお、Describe とは、説明的な関係にあり、元文書と先文書が共存しないと文脈的には理解できないほど強い結びつきを持っていることを意味する Link である。例えば、添付ファイルとそれを説明するメッセージがこの関係に当てはまる。Follow とは、ある文書を受けてそれに補足 / 返答 / 反論するなど、文脈的な情報を追加していくこと全般を意味する Link である。Revise とは、ある文書を改訂して置き換えている関係を意味する Link である。Sum とは、ある文書をまとめたり、要約したりする関係を意味する Link である。Consult とは、ある文書を参考にしたという関係を意味する Link である。

40

【 0 0 5 7 】

続いて、図 6 乃至図 1 5 を参照して、図 4 のステップ S 4 におけるキーワード候補生成部 6 0 によるキーワード候補生成処理について説明する。

【 0 0 5 8 】

図 6 は本実施形態のキーワード候補生成部 6 0 によるキーワード候補生成処理の一例を示すフローチャートである。

【 0 0 5 9 】

キーワード候補生成部 6 0 が設定情報を取得すると、注目文書検索部 6 1 が、この設定情報に基づいて文書を抽出する（ステップ S 1 1）。具体的には、注目文書検索部 6 1 は

50

、設定情報に登録されている検索キーワード102及び検索対象の期間106、検索対象のタスク名と合致する文書をタスクDB14と文書DB15から抽出する。条件に該当し、抽出された文書を注目文書とする。そして注目文書検索部61は、抽出した注目文書データはカテゴリ判定部62に送信する。なお、設定情報の取得は、設定情報DB13に登録された設定情報を取得してもよいし、入力部31から入力された設定情報を直接受信してもよい。

【0060】

続いて、カテゴリ判定部62は、カテゴリDB16を参照して注目文書検索部61から受信した注目文書データのカテゴリを解析する(ステップS12)。具体的には、カテゴリ判定部62は、カテゴリDB16からカテゴリ情報を取得する。そして、取得したカテゴリ情報と注目文書検索部61から受信し注目文書データの一覧とに基づいて、カテゴリ判定部62は、検索結果の注目文書ごとに、どのカテゴリに属しているかを解析する。カテゴリ判定部62による解析結果はキーワード候補抽出部62に送信される。なお、各注目文書がどのカテゴリに属しているかという情報をカテゴリ解析結果と呼ぶものとする。

10

【0061】

図7にカテゴリ解析結果の概念図を示す。図7に示すように、本実施形態の情報検索システムにおいては、カテゴリDB16に「HTML」161、「ネットワーク」162、および「無線LAN」163のカテゴリが存在する。HTML161のカテゴリには、文書データ151-1乃至151-5が属している。そのうち、当該カテゴリに属すると判定される文書データは、斜線の文書である注目文書データ151-1乃至151-3である。また、ネットワーク162のカテゴリには、文書データ152-1乃至152-3が属している。そのうち、当該カテゴリに属すると判定される文書データは、ドットの文書である注目文書データ152-1である。また、無線LAN163のカテゴリには、文書データ153-1乃至153-6が属している。そのうち、当該カテゴリに属すると判定される文書データはない。

20

【0062】

キーワード候補抽出部63が解析結果を受信すると、キーワード候補抽出部63によって、ステップS13乃至ステップS17のキーワード候補抽出処理が行われる。

【0063】

まず、キーワード候補抽出部63のカテゴリ名取得部66がキーワード候補の抽出を行う(ステップS13)。次に、キーワード候補抽出部63の新カテゴリ名取得部67がキーワード候補の抽出を行う(ステップS14)。続いて、キーワード候補抽出部63の出現頻度解析部64がキーワード候補の抽出を行う(ステップS15)。続いて、キーワード候補抽出部63の期間限定出現頻度解析部65がキーワード候補の抽出を行う(ステップS16)。最後に、キーワード候補抽出部63のKL連携解析部68がキーワード候補の抽出を行う(ステップS17)。

30

【0064】

キーワード候補抽出処理が終了すると、キーワード候補抽出部63は抽出されたキーワード候補を全てクライアント端末20に送信する(ステップS18)。これにより、キーワード候補生成処理は終了する。

40

【0065】

ここで、図8乃至図15を参照して、図6のステップS12乃至ステップS17のそれぞれの処理を詳しく説明する。

【0066】

まず、図8を参照して、図6のステップS13におけるカテゴリ名取得部66によるキーワード候補抽出処理について説明する。図8はカテゴリ名取得部66によるキーワード候補抽出処理の一例を示すフローチャートである。

【0067】

まず、カテゴリ名取得部66は、カテゴリ判定部62によるカテゴリ解析結果を受信する(ステップS131)。受信したカテゴリ解析結果に基づいて、カテゴリ名取得部66

50

は、検索結果の注目文書群のカテゴリごとの個数を集計する（ステップ S 1 3 2）。

【 0 0 6 8 】

集計した結果、注目文書が特定数以上含まれている場合（ステップ S 1 3 3 が Yes）、該当するカテゴリのカテゴリ名をキーワード候補と判定する（ステップ S 1 3 4）。ステップ S 1 3 4 でキーワード候補と判定された全てのカテゴリ名をキーワード候補抽出部 6 3 に送信し（ステップ S 1 3 5）、処理を終了する。

【 0 0 6 9 】

本実施形態では、注目文書が 3 以上含まれているカテゴリのカテゴリ名を抽出するとする。すなわち、「HTML」が抽出される。なお、集計した結果、注目文書が特定数以上含まれていない場合（ステップ S 1 3 3 が No）、キーワードの候補は抽出せず、処理を終了する。

10

【 0 0 7 0 】

本実施形態では、カテゴリ名をキーワード候補とすることで、検索キーワードと関連付いたキーワードの抽出を可能とする。

【 0 0 7 1 】

図 9 を参照して、図 6 のステップ S 1 4 における新カテゴリ名取得部 6 7 によるキーワード候補抽出処理について説明する。図 9 は新カテゴリ名取得部 6 7 によるキーワード候補抽出処理の一例を示すフローチャートである。

【 0 0 7 2 】

まず、新カテゴリ名取得部 6 7 は、カテゴリ判定部 6 2 によるカテゴリ解析結果を受信する（ステップ S 1 4 1）。受信したカテゴリ解析結果に基づいて、新カテゴリ名取得部 6 7 は、検索結果の注目文書群のカテゴリごとの個数を集計する（ステップ S 1 4 2）。

20

【 0 0 7 3 】

集計した結果、注目文書が特定数以上含まれているカテゴリが存在する場合（ステップ S 1 4 3 が Yes）、新カテゴリ名取得部 6 7 は、当該カテゴリに含まれる文書データのクラスタリング分析を行い、新たなカテゴリ（以下、新カテゴリという）を生成する（ステップ S 1 4 4）。

【 0 0 7 4 】

新カテゴリが生成されると、新カテゴリ名取得部 6 7 は、新カテゴリの中に、注目文書データが含まれている新カテゴリがあるか否かを判定する（ステップ S 1 4 5）。注目文書データが含まれている新カテゴリがある場合（ステップ S 1 4 5 が Yes）対象の新カテゴリ名（以下、新カテゴリ名という）をキーワード候補と判定し（ステップ S 1 4 6）、当該キーワード候補をキーワード候補抽出部 6 3 に送信して（ステップ S 1 4 7）、処理を終了する。

30

【 0 0 7 5 】

なお、集計した結果、注目文書が特定数以上含まれているカテゴリが存在しない場合（ステップ S 1 4 3 が No）、キーワードの候補は抽出せず、処理を終了する。また、注目文書データが含まれている新カテゴリがない場合（ステップ S 1 4 5 が No）も、キーワードの候補は抽出せず、処理を終了する。

【 0 0 7 6 】

図 10 を参照して、本実施形態の情報検索システムにおける新カテゴリ名取得部 6 7 によるキーワード候補抽出処理の一例を具体的に説明する。なお、本実施形態では、あらかじめ定められたカテゴリ毎に属する文書数の閾値は 3 である。

40

【 0 0 7 7 】

したがって、図 10 に示すように、HTML 1 6 1 に属する文書データに関して処理が行われる。すなわち、HTML 1 6 1 に属する文書データにクラスタリングが行われ、カテゴリが生成される。ここでは、新カテゴリとして、タグ 1 6 4 とハイパーテキスト 1 6 5 とが生成されている。このタグ 1 6 4 には注目文書データが、属しているため、本実施形態ではタグ 1 6 4 がキーワード候補として抽出される。

【 0 0 7 8 】

50

次に図 1 1 を参照して、図 6 のステップ S 1 5 における出現頻度解析部 6 4 によるキーワード候補抽出処理について説明する。図 1 1 は出現頻度解析部 6 4 によるキーワード候補抽出処理の一例を示すフローチャートである。

【 0 0 7 9 】

まず、出現頻度解析部 6 4 は、カテゴリ判定部 6 2 によるカテゴリ解析結果を受信する（ステップ S 1 5 1）。受信したカテゴリ解析結果に基づいて、出現頻度解析部 6 4 は、検索結果の注目文書が属しているカテゴリ毎に、当該カテゴリ内の文書データにおいて出現する頻度の高い単語を抽出する。

【 0 0 8 0 】

すなわち、出現頻度解析部 6 4 は、検索結果の注目文書が属しているカテゴリ毎に、当該カテゴリ内の文書データに含まれる各単語のうち、検索キーワード 1 0 2 に設定されている単語以外の単語についての出現頻度を集計する（ステップ S 1 5 2）。 10

【 0 0 8 1 】

集計した結果、出現頻度が特定数以上の単語が存在する場合（ステップ S 1 5 3 が Yes）、これらの単語をキーワード候補と判定し（ステップ S 1 5 4）、当該キーワード候補をキーワード候補抽出部 6 3 に送信して（ステップ S 1 5 5）、処理を終了する。

【 0 0 8 2 】

集計した結果、出現頻度が特定数以上の単語が存在しない場合（ステップ S 1 5 3 が No）、出現頻度解析部 6 4 はキーワードの候補は抽出せず、処理を終了する。

【 0 0 8 3 】

図 1 2 に、出現頻度解析部 6 4 による HTML 1 6 1 カテゴリからのキーワード候補抽出処理の概念図を示す。図 1 2 に示すように、出現頻度解析部 6 4 は HTML 1 6 1 カテゴリに属する文書データのタイトルに含まれる単語を抽出する。そのうち、出現数があらかじめ定めた数値より多いものであって、検索キーワード 1 0 2 である「インターネット」を除く単語をキーワード候補とする。すなわち、網掛けの「グリッドコンピュータ」という単語がキーワード候補として抽出される。 20

【 0 0 8 4 】

検索キーワード 1 0 2 が属しているカテゴリにおける出現頻度の高い単語は、ユーザが所望する情報との関連性が高いと考えられる。このため、これらの単語をキーワード候補としてユーザに提示することで、ユーザが所望する情報を取得する精度を上げることができる。 30

【 0 0 8 5 】

図 1 3 を参照して、図 6 のステップ S 1 6 における期間限定出現頻度解析部 6 5 によるキーワード候補抽出処理について説明する。図 1 3 は期間限定出現頻度解析部 6 5 によるキーワード候補抽出処理の一例を示すフローチャートである。

【 0 0 8 6 】

まず、期間限定出現頻度解析部 6 5 は、カテゴリ判定部 6 2 によるカテゴリ解析結果を受信する（ステップ S 1 6 1）。受信したカテゴリ解析結果に基づいて、期間限定出現頻度解析部 6 5 は、検索結果の注目文書が属しているカテゴリ毎に、当該カテゴリ内の文書データにおいて、設定情報の期間 1 0 6 によって設定された期間内に生成された、すなわち、文書 DB 1 5 に登録された文書を抽出する（ステップ S 1 6 2）。すなわち、期間限定出現頻度解析部 6 5 は、対象のカテゴリ内の文書データを文書の登録日時によって絞り込む。 40

【 0 0 8 7 】

ステップ S 1 6 2 において抽出された文書データにおいて、期間限定出現頻度解析部 6 5 は、検索キーワード 1 0 2 に設定されている単語以外の単語において出現する頻度の高い単語を抽出する。すなわち、期間限定出現頻度解析部 6 5 は、検索結果の注目文書が属しているカテゴリ毎に、当該カテゴリ内の文書データに含まれる各単語の出現頻度を集計する（ステップ S 1 6 3）。

【 0 0 8 8 】

集計した結果、出現頻度が特定数以上の単語が存在する場合（ステップ S 1 6 4 が Yes）、これらの単語をキーワード候補と判定し（ステップ S 1 6 5）、当該キーワード候補をキーワード候補抽出部 6 3 に送信して（ステップ S 1 6 6）、処理を終了する。

【 0 0 8 9 】

集計した結果、出現頻度が特定数以上の単語が存在しない場合（ステップ S 1 6 4 が No）、期間限定出現頻度解析部 6 5 はキーワードの候補は抽出せず、処理を終了する。

【 0 0 9 0 】

検索対象となる文書データを指定した期間 1 0 2 で絞り込むことで、最新の文書を対象にして検索キーワードの候補を抽出することが可能となる。したがって、これらのキーワード候補をユーザに提示することで、ユーザが所望する情報を取得する精度を上げることができる。

10

【 0 0 9 1 】

図 1 4 を参照して、図 6 のステップ S 1 7 における K L 連携解析部 6 8 によるキーワード候補抽出処理について説明する。図 1 4 は K L 連携解析部 6 8 によるキーワード候補抽出処理の一例を示すフローチャートである。

【 0 0 9 2 】

まず、K L 連携解析部 6 8 は、カテゴリ判定部 6 2 によるカテゴリ解析結果を受信する（ステップ S 1 7 1）。受信したカテゴリ解析結果に基づいて、K L 連携解析部 6 8 は、検索結果の注目文書が属しているカテゴリ内の文書データにおいて、当該注目文書とナレッジリンクにより関連付いている文書を抽出する（ステップ S 1 7 2）。このナレッジリンクに関する情報は文書 DB 1 5 から抽出される。

20

【 0 0 9 3 】

ステップ S 1 6 2 において抽出された文書データにおいて、K L 連携解析部 6 8 は、検索キーワード 1 0 2 に設定されている単語以外の出現する頻度の高い単語を抽出する。すなわち、K L 連携解析部 6 8 は、抽出結果の文書データが属しているカテゴリ毎に、当該カテゴリ内の文書データに含まれる検索キーワード 1 0 2 に設定されている単語以外の単語の各出現頻度を集計する（ステップ S 1 7 3）。

【 0 0 9 4 】

集計した結果、出現頻度が特定数以上の単語が存在する場合（ステップ S 1 7 4 が Yes）、これらの単語をキーワード候補と判定し（ステップ S 1 7 5）、当該キーワード候補をキーワード候補抽出部 6 3 に送信して（ステップ S 1 7 6）、処理を終了する。

30

【 0 0 9 5 】

集計した結果、出現頻度が特定数以上の単語が存在しない場合（ステップ S 1 7 4 が No）、K L 連携解析部 6 8 はキーワードの候補は抽出せず、処理を終了する。

【 0 0 9 6 】

図 1 5 に、K L 連携解析部 6 8 による HTML 1 6 1 カテゴリからのキーワード候補抽出処理の概念図を示す。図 1 5 に示すように、K L 連携解析部 6 8 は、HTML 1 6 1 に属する文書データに、K L によって関連付けられた文書データを抽出する（点線の文書）。K L 連携解析部 6 8 は、これらの抽出された文書データのタイトルに含まれる単語であって、あらかじめ定めた数値以上含まれている単語をキーワード候補として抽出する。この際、検索キーワード 1 0 2 に入力された単語は除外される。

40

【 0 0 9 7 】

上述したように、K L 連携解析部 6 8 は、ユーザが指定した検索キーワード 1 0 2 を含む文書データと関連する文書データが属するカテゴリにおいて、特定数以上の出現頻度のある単語をキーワード候補として抽出する。すなわち、本実施形態は、ユーザが検索した文書と関連する文書における出現頻度の高い単語をキーワード候補としてユーザに提示することで、ユーザが所望する情報を取得する精度を上げることが可能とする。

【 0 0 9 8 】

図 1 6 に本実施形態の情報検索システムにおけるキーワード候補生成処理後に表示部 3 2 に表示される画面の一例を示す。

50

【0099】

図16に示すように、本実施形態の情報検索システムは、キーワード候補生成処理後、検索条件設定画面100の推奨キーワード103に、キーワード候補抽出部63によって抽出されたキーワード候補を表示する。本実施形態においては、「HTML」と「グリッドコンピューティング」とが表示される。

【0100】

すなわち、本実施形態の情報検索システムは、このようにユーザが入力した検索キーワードに基づいて、検索キーワードの候補を提示することによって、ユーザが指定した検索キーワードが廃れた場合など、ユーザが入力した検索キーワード以外のよりユーザの注目している文書を検索できるキーワードをユーザに提示することができる。これにより、ユーザ新たに検索キーワードを設定する必要をなくすことを可能とする。また、ユーザが知らない最新の検索キーワードを取得することを可能とする。

10

【0101】

故に、本実施形態の情報検索システムは、検索において必要とされる最新の情報を取得することが可能となるため、業務効率を向上させることができる。また、ユーザが所望する情報を検索する精度を向上する検索キーワードをユーザに提示することを可能とする。

【0102】

なお、本実施形態の情報検索システムは、検索キーワードと関連付いている単語を探すような場合に、単語間の関連性に関する情報を保持する関連辞書の作成やメンテナンスなどの作業を行う必要が無い。また本実施形態の情報検索システムは、ユーザが検索で用いたキーワードを元にどの文書を参照したかを解析するためのログ情報の取得やログ情報の解析などの処理をする必要も無い。

20

【0103】

なお、ここでは本発明の実施形態を説明したが、この実施形態は例として提示したものであり、発明の範囲を限定することは意図していない。この新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。この実施形態やその変形は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

【0104】

例えば、キーワード候補抽出処理における、単語の出現頻度が高いかを判定する際に用いられる特定数は、本実施形態の情報検索システム構築時に設定してもいいし、サーバの管理者によって設定できるようにしてもよい。

30

また、例えば、本実施形態の情報検索システムにおける文書検索は、ユーザの指定したタイミングで行われるようにしても良い。

【0105】

また、キーワード候補抽出処理は、キーワード候補抽出部63が有する各部が順次行ってもよいし、同時に行っても良い。なお、キーワード候補抽出部63が有する各部がキーワード候補抽出処理を行う順番は特に限定されない。また、キーワード候補抽出処理はすべてを行わず、いくつかの処理の組み合わせでも良い。

【0106】

また、キーワード候補抽出部63によるキーワード候補抽出処理において、期間限定出現頻度解析部65における単語の出現頻度の集計は行わず、出現頻度解析部64による単語の出現頻度の集計データを用いても良い。

40

【0107】

また、キーワード候補抽出部63によるキーワード候補抽出処理における単語の出現頻度の集計は、出現頻度解析部64、期間限定出現頻度解析部65、及びKL連携解析部67それぞれで行わず、例えば、出現頻度解析部64などのどれか一つの部による単語の出現頻度の集計データをその他の部においても用いるようにしても良い。

【0108】

同様に、キーワード候補抽出部63によるキーワード候補抽出処理における各カテゴリ

50

に含まれている注目文書の集計は、カテゴリ名取得部 66、新カテゴリ名取得部 67 それぞれで行わず、例えば、カテゴリ名取得部 66 による計データを新カテゴリ名取得部 67 においても用いるようにしても良い。

【0109】

また、検索条件はその検索の周期と対応させて、複数登録することが可能である。

【0110】

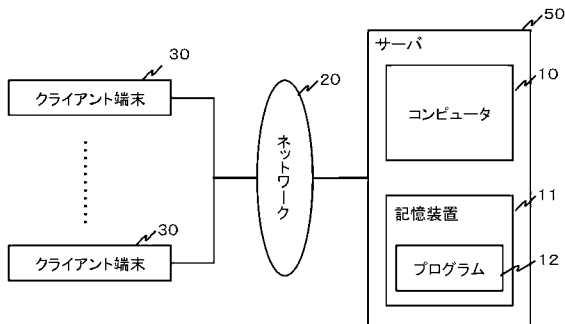
また、設定情報には、文書データの作成者や、文書データのラベルなどを設定してもよい。

【符号の説明】

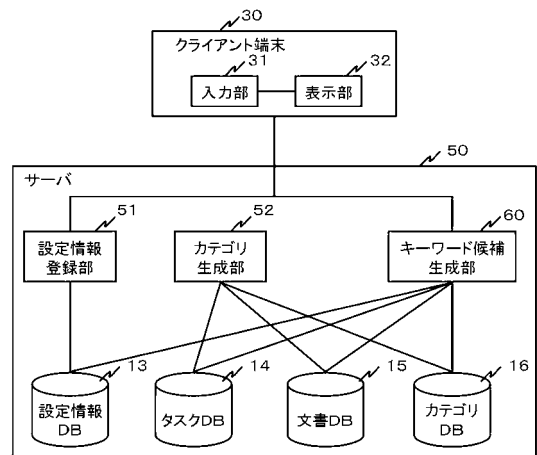
【0111】

10...コンピュータ、11...記憶装置、12...プログラム、13...設定情報DB、14...タスクDB、15...文書DB、16...カテゴリDB、20...ネットワーク、30...クライアントコンピュータ、31...入力部、32...表示部、50...サーバ、51...設定情報登録部、52...カテゴリ生成部、60...キーワード候補生成部、61...注目文書検索部、62...カテゴリ判定部、63...キーワード候補抽出部、64...出現頻度解析部、65...期間限定出現頻度解析部、66...カテゴリ名取得部、67...新カテゴリ名取得部、68...ナレッジリンク連携取得部

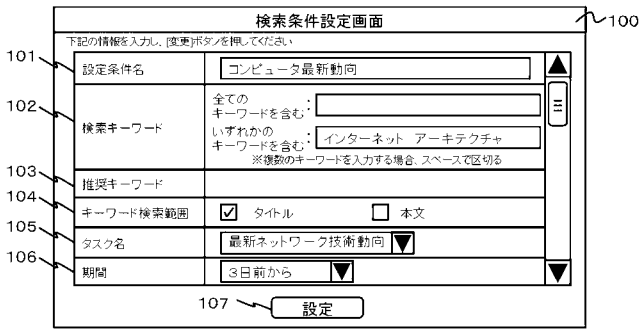
【図1】



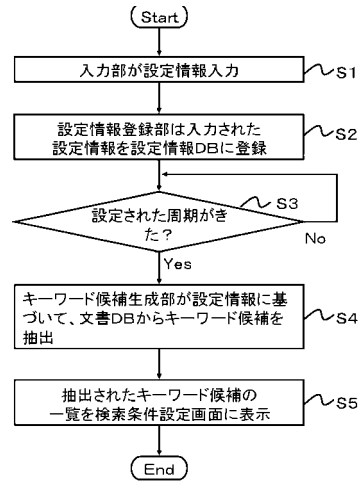
【図2】



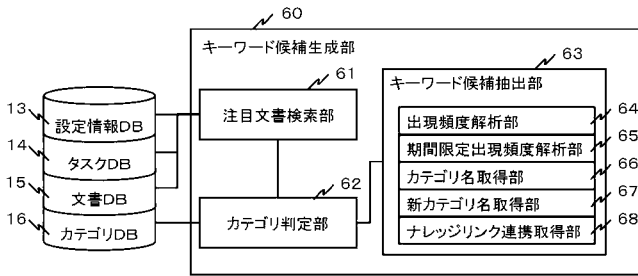
【 図 3 】



【 図 4 】



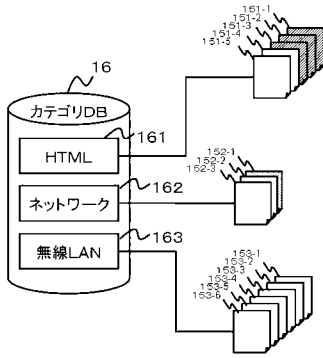
【 図 5 】



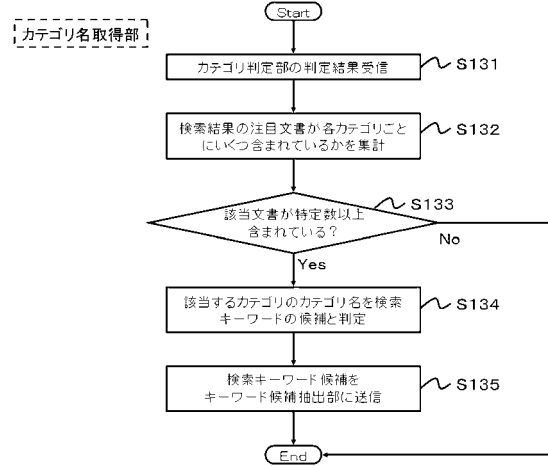
【 図 6 】



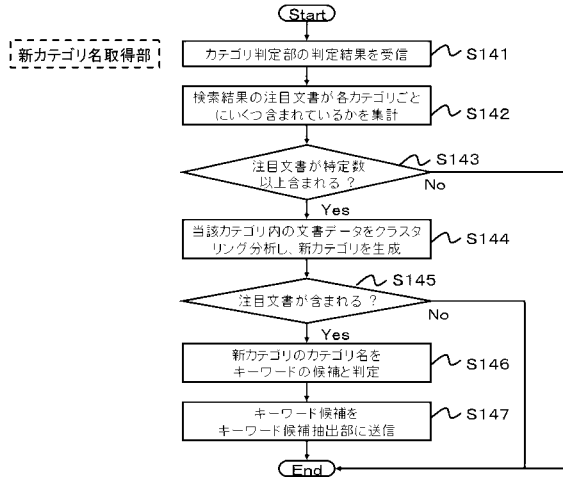
【 図 7 】



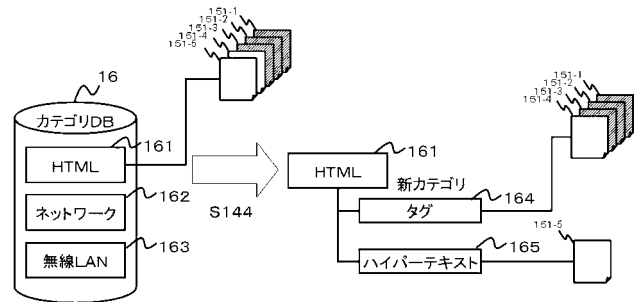
【 図 8 】



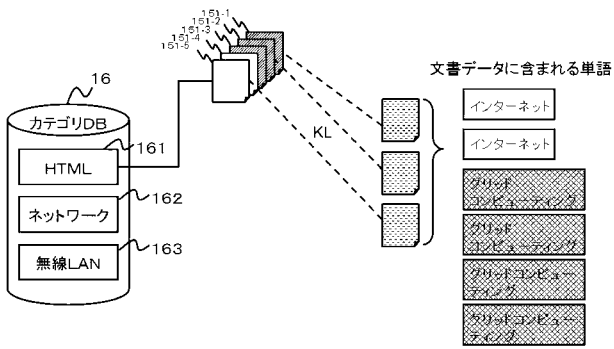
【 図 9 】



【 図 10 】



【 図 1 5 】



【 図 1 6 】

Figure 16 shows a '検索条件設定画面' (Search Condition Setting Screen) with various input fields and controls. The screen includes a title bar (100) and a subtitle '下記の情報を入力し、変更ボタンを押してください' (Please enter the following information and press the change button). The fields are numbered as follows:

- 101: 設定条件名 (Setting Condition Name) - コンピュータ最新動向 (Computer Latest Trends)
- 102: 検索キーワード (Search Keyword) - 全てのキーワードを含む (Include all keywords) and いずれかのキーワードを含む (Include any keyword). Example: インターネット アーキテクチャ (Internet Architecture). Note: ※複数のキーワードを入力する場合、スペースで区切る (When entering multiple keywords, separate with spaces).
- 103: 追加キーワード (Additional Keyword) - HTML グリッドコンピューティング (HTML Grid Computing)
- 104: キーワード検索範囲 (Keyword Search Range) - タイトル (Title) 本文 (Body)
- 105: タスク名 (Task Name) - 最新ネットワーク技術動向 (Latest Network Technology Trends)
- 106: 期間 (Period) - 3日前から (From 3 days ago)
- 107: 設定 (Settings) button

フロントページの続き

(72)発明者 佐々木 淳哉

東京都港区芝浦一丁目1番1号 東芝ソリューション株式会社内

Fターム(参考) 5B075 ND03 QM10 UU40