

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5494066号
(P5494066)

(45) 発行日 平成26年5月14日(2014.5.14)

(24) 登録日 平成26年3月14日(2014.3.14)

(51) Int.Cl. F I
G06F 17/30 (2006.01)
 G06F 17/30 4 1 4 Z
 G06F 17/30 1 7 0 A
 G06F 17/30 4 1 4 B

請求項の数 8 (全 13 頁)

(21) 出願番号	特願2010-61451 (P2010-61451)	(73) 特許権者	000005223 富士通株式会社
(22) 出願日	平成22年3月17日(2010.3.17)		神奈川県川崎市中原区上小田中4丁目1番1号
(65) 公開番号	特開2011-197809 (P2011-197809A)	(74) 代理人	100089118 弁理士 酒井 宏明
(43) 公開日	平成23年10月6日(2011.10.6)	(72) 発明者	大堀 順也 福岡県福岡市早良区百道浜二丁目2番1号 株式会社富士通九州システムズ内
審査請求日	平成25年1月8日(2013.1.8)	審査官	野崎 大進

最終頁に続く

(54) 【発明の名称】 検索装置、検索方法および検索プログラム

(57) 【特許請求の範囲】

【請求項1】

意味を持つ単語毎に文字列を区切る単語区切方式に基づいて区切られ、文書データに関連付けられた単語インデックスと、文字毎に文字列を区切る文字区切方式に基づいて区切られ、文書データに関連付けられた文字インデックスと、記号を含む所定の文字列の特徴を定義したパターンファイルとを記憶する記憶部と、

検索文字列を受け付け、前記検索文字列と前記パターンファイルとを基にして、前記パターンファイルに定義された特徴が、前記検索文字列と一致する場合には、前記単語インデックスを用いて文書データの検索を行うと判定し、前記パターンファイルに定義された特徴が、前記検索文字列と一致しない場合には、前記文字インデックスを用いて文書データの検索を行うと判定する判定部と、

前記判定部の判定結果に基づいて、前記単語インデックスまたは前記文字インデックスを用いて文書データの検索を実行する検索部と

を備えたことを特徴とする検索装置。

【請求項2】

前記パターンファイルは前記所定の文字列に含まれる文字の形式を示す情報を含み、前記判定部は、受け付けた前記検索文字列に含まれる文字の形式と前記パターンファイルに示される文字の形式とが一致するか否かに基づき、前記単語インデックスを用いて文書データの検索を行うのか、前記文字インデックスを用いて文書データの検索を行うのかを判定する

ことを特徴とする請求項 1 に記載の検索装置。

【請求項 3】

前記検索文字列を、複数の部分文字列に分割する文字列分割部を更に有し、前記判定部は、部分文字列毎に、前記パターンファイルに定義された特徴が、前記部分文字列と一致する場合には、前記単語インデックスを用いて文書データの検索を行うと判定し、前記パターンファイルに定義された特徴が、前記部分文字列と一致しない場合には、前記文字インデックスを用いて文書データの検索を行うと判定することを特徴とする請求項 1 または 2 に記載の検索装置。

【請求項 4】

意味を持つ単語毎に文字列を区切る単語区切方式に基づいて区切られ、文書データに関連付けられた単語インデックスと、文字毎に文字列を区切る文字区切方式に基づいて区切られ、文書データに関連付けられた文字インデックスと、記号を含む所定の文字列の特徴を定義したパターンファイルとを記憶する記憶装置を有する検索装置が、

検索文字列を受け付け、前記検索文字列と前記パターンファイルとを基にして、前記パターンファイルに定義された特徴が、前記検索文字列と一致する場合には、前記単語インデックスを用いて文書データの検索を行うと判定し、前記パターンファイルに定義された特徴が、前記検索文字列と一致しない場合には、前記文字インデックスを用いて文書データの検索を行うと判定する判定ステップと、

前記判定ステップの判定結果に基づいて、前記単語インデックスまたは前記文字インデックスを用いて文書データの検索を実行する検索ステップと

を含むことを特徴とする検索方法。

【請求項 5】

意味を持つ単語毎に文字列を区切る単語区切方式に基づいて区切られ、文書データに関連付けられた単語インデックスと、文字毎に文字列を区切る文字区切方式に基づいて区切られ、文書データに関連付けられた文字インデックスと、記号を含む所定の文字列の特徴を定義したパターンファイルとを記憶する記憶装置を有するコンピュータに、

検索文字列を受け付け、前記検索文字列と前記パターンファイルとを基にして、前記パターンファイルに定義された特徴が、前記検索文字列と一致する場合には、前記単語インデックスを用いて文書データの検索を行うと判定し、前記パターンファイルに定義された特徴が、前記検索文字列と一致しない場合には、前記文字インデックスを用いて文書データの検索を行うと判定する判定手順と、

前記判定ステップの判定結果に基づいて、前記単語インデックスまたは前記文字インデックスを用いて文書データの検索を実行する検索手順と

を実行させることを特徴とする検索プログラム。

【請求項 6】

第 1 の区切方式に基づいて区切られ、文書データに関連付けられた第 1 のインデックスと、第 2 の区切方式に基づいて区切られ、文書データに関連付けられた第 2 のインデックスと、文字列に含まれる文字の形式を示すパターンファイルとを記憶する記憶部と、

検索文字列を受け付け、前記検索文字列に含まれる文字の形式と前記パターンファイルに示される文字の形式とが一致するか否かに基づき、前記第 1 のインデックスを用いて文書データの検索を行うのか、前記第 2 のインデックスを用いて文書データの検索を行うのかを判定する判定部と、

前記判定部の判定結果に基づいて、前記第 1 のインデックスまたは前記第 2 のインデックスを用いて文書データの検索を実行する検索部と

を備えたことを特徴とする検索装置。

【請求項 7】

第 1 の区切方式に基づいて区切られ、文書データに関連付けられた第 1 のインデックスと、第 2 の区切方式に基づいて区切られ、文書データに関連付けられた第 2 のインデックスと、文字列に含まれる文字の形式を示すパターンファイルとを記憶する記憶装置を有する検索装置が、

10

20

30

40

50

検索文字列を受け付け、前記検索文字列に含まれる文字の形式と前記パターンファイルに示される文字の形式とが一致するか否かに基づき、前記第1のインデックスを用いて文書データの検索を行うのか、前記第2のインデックスを用いて文書データの検索を行うのかを判定する判定ステップと、

前記判定ステップの判定結果に基づいて、前記第1のインデックスまたは前記第2のインデックスを用いて文書データの検索を実行する検索ステップと
を含むことを特徴とする検索方法。

【請求項8】

第1の区切方式に基づいて区切られ、文書データに関連付けられた第1のインデックスと、第2の区切方式に基づいて区切られ、文書データに関連付けられた第2のインデックスと、文字列に含まれる文字の形式を示すパターンファイルとを記憶する記憶装置を有するコンピュータに、

10

検索文字列を受け付け、前記検索文字列に含まれる文字の形式と前記パターンファイルに示される文字の形式とが一致するか否かに基づき、前記第1のインデックスを用いて文書データの検索を行うのか、前記第2のインデックスを用いて文書データの検索を行うのかを判定する判定手順と、

前記判定手順の判定結果に基づいて、前記第1のインデックスまたは前記第2のインデックスを用いて文書データの検索を実行する検索手順と

を実行させることを特徴とする検索プログラム。

【発明の詳細な説明】

20

【技術分野】

【0001】

本発明は、検索装置等に関する。

【背景技術】

【0002】

複数の文書データから特定の文字列を検索する全文検索が知られている。この全文検索では、転置インデックスが用いられる。転置インデックスは、文字データに含まれる単語の位置情報等を格納する索引に対応する。転置インデックスを作成する方式には、大きく分けて文字区切方式と、単語区切方式とがある。

【0003】

30

文字区切方式では、単語の意味を考えずに、文字単位で転置インデックスを作成するものである。文字区切方式で作成した転置インデックスを文字インデックスと表記する。文字インデックスを用いれば、完全な部分一致検索が可能である。しかし、検索キーワードと文字インデックスとを一文字ずつ比較する必要があり、検索時間を多く要してしまうという欠点がある。

【0004】

単語区切方式では、意味のある単語単位で転置インデックスを作成するものである。単語区切方式で作成した転置インデックスを単語インデックスと表記する。単語インデックスを用いれば、検索キーワードを単語毎に比較するので、文字インデックスを利用する場合と比較して、検索時間を短縮することが可能となる。しかし、単語の区切り方によって、検索漏れが発生する場合がある。

40

【0005】

このように、文字区切方式および単語区切方式には、それぞれ長所、短所があるため、いかにして文字区切方式と単語区切方式とを使い分けるのが重要になる。例えば、文字区切方式および単語区切方式を用いた従来技術として、検索キーワードの長さに応じて、文字インデックスと単語インデックスとを自動選択するという技術が開示されている。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開平10-307835号公報

50

【特許文献2】特開2001-34623号公報

【特許文献3】特開2008-77673号公報

【発明の概要】

【発明が解決しようとする課題】

【0007】

しかしながら、全文検索を行う文書データによっては、検索キーワードの長さが同じ場合でも、文字インデックスと単語インデックスとを使い分けた方が効率的な全文検索を行える場合がある。

【0008】

例えば、バイオデータベースに記憶される文書データには、文書に加えて他のデータベースへリンクするためのID (Identification) が含まれている。一般的に、ID等の記号を有さない文書データに対しては、単語インデックスが有効であり、記号を有する文書データに対しては、文字インデックスが有効である。

【0009】

ここで、「1.1.1.1ANDsuppressor」という検索式が与えられた場合を例にして説明する。かかる検索式に対して、文字インデックスを用いて全文検索を試みる場合には、「1.1.1.1」というIDを含み、かつ、「suppressor」という単語を含む文書データのみを検索することが好ましい。しかし、上記検索式に対して、文字インデックスを用いて全文検索を試みると、実際には、「1.1.1.11」、「1.1.1.12」等のIDを含む文書データもヒット

【0010】

これに対して、上記検索式に対して、単語インデックスを用いて全文検索を試みると、「1.1.1.1」のIDを含む文書データのみを検索することが可能である。しかし、「suppressors」と「suppressor」とは完全に一致していないので、「suppressors」を含む文書データを検索することが出来なくなってしまう。

【0011】

開示の技術は、上記に鑑みてなされたものであって、文書データの特性によらず、効率よく全文検索を実行することができる検索装置、検索方法および検索プログラムを提供することを目的とする。

【課題を解決するための手段】

【0012】

本願の開示する検索装置は、一つの態様において、第1の区切方式に基づいて区切られ、文書データに関連付けられた第1のインデックスと、第2の区切方式に基づいて区切られ、文書データに関連付けられた第2のインデックスと、所定の文字の特徴を定義したパターンファイルを記憶する記憶部と、検索キーワードを受け付け、前記検索キーワードと前記パターンファイルとを基にして、前記第1のインデックスを用いて文書データの検索を行うのか、前記第2のインデックスを用いて文書データの検索を行うのかを判定する判定部と、前記判定部の判定結果に基づいて、前記第1のインデックスまたは前記第2のインデックスを用いて文書データの検索を実行する検索部とを備えたことを要件とする。

【発明の効果】

【0013】

本願の開示する検索装置の一つの態様によれば、文書データの特性によらず、効率よく全文検索を実行することができるという効果を奏する。

【図面の簡単な説明】

【0014】

【図1】図1は、本実施例1にかかる検索装置の構成を示す図である。

【図2】図2は、本実施例2にかかるシステムを示す図である。

【図3】図3は、本実施例2にかかる検索装置の構成を示す図である。

【図4】図4は、パターンファイルのデータ構造を示す図である。

【図5】図5は、本実施例2にかかる検索装置の処理手順を示すフローチャートである。

10

20

30

40

50

【図6】図6は、実施例にかかる検索装置を構成するコンピュータのハードウェア構成を示す図である。

【発明を実施するための形態】

【0015】

以下に、本願の開示する検索装置、検索方法および検索プログラムの実施例を図面に基づいて詳細に説明する。なお、この実施例によりこの発明が限定されるものではない。

【実施例1】

【0016】

図1は、本実施例1にかかる検索装置100の構成を示す図である。図1に示すように、この検索装置100は、記憶部110、判定部120、検索部130を有する。

10

【0017】

記憶部110は、パターンファイル110a、第1のインデックス110b、第2のインデックス110cを記憶する。パターンファイル110aは、所定の文字の特徴を定義したデータである。第1のインデックス110bは、第1の区切方式に基づいて区切られ、文書データに関連付けられたデータである。第2のインデックス110cは、第2の区切方式に基づいて区切られ、文書データに関連付けられたデータである。

【0018】

判定部120は、検索キーワードを受け付け、検索キーワードとパターンファイル110aとを基にして、第1のインデックス110bを用いて検索を行うのか、第2のインデックス110cを用いて検索を行うのかを判定する。

20

【0019】

検索部130は、判定部120の判定結果に基づいて、第1のインデックス110bまたは第2のインデックス110cを用いて文書データの検索を実行する。

【0020】

上記の検索装置100は、パターンファイル110aを用いて、第1のインデックス110bを用いた検索を行うのか、第2のインデックス110cを用いた検索を行うのかを判定している。このため、検索キーワードの特徴に合わせて最適なインデックスを選択することができるので、文書データの特性によらず、効率よく全文検索を実行することができる。

【実施例2】

30

【0021】

次に、本実施例2にかかるシステムの一例について説明する。図2は、本実施例2にかかるシステムを示す図である。図2に示すように、このシステムは、利用者端末60、検索装置200を有する。利用者端末60と検索装置200は、ネットワーク50を介して接続される。

【0022】

利用者端末60は、検索装置200に検索キーワードを送信し、検索キーワードに対する検索結果を検索装置200から受信する装置である。

【0023】

検索装置200は、文書データの全文検索を行う装置である。図3は、本実施例2にかかる検索装置200の構成を示す図である。図3に示すように、この検索装置200は、記憶部210、インデクシング処理部220、入力受付部230、検索式解析処理部240、スコアリング処理部250、検索結果出力部260を有する。

40

【0024】

記憶部210は、パターンファイル210a、文書データ群210b、単語インデックス210c、文字インデックス210dを記憶する。

【0025】

パターンファイル210aは、所定の文字の特徴を定義したデータである。図4は、パターンファイル210aのデータ構造を示す図である。図4に示すように、このパターンファイルは、Noとパターンとを有する。Noは、各パターンを識別するものである。パ

50

ターンは、所定の文字の特徴を正規表現で示したものである。ここで、文字には、一般的な文字のほかに、数字や記号等も含まれるものとする。

【0026】

ここで、パターンの記載方法の一例について説明する。パターン中の[]は、[と]の中に書かれたいずれかの一文字に一致する文字、数字、記号を意味する。例えば、[0-9]は、1桁の数字を意味する。パターン中の{n, m}は、直前の文字がn回からm回まで繰り返されることを意味する。例えば、[0-9]{1, 3}は、1桁、2桁、3桁の数字を意味する。

【0027】

また、パターン中の+は、直前の文字が1回以上繰り返されることを意味する。例えば、[0-9]+は、数字からなる文字列を意味する。パターン中の*は、直前の文字が0回以上繰り返されることを意味する。例えば、[0-9]*は、空文字または数字からなる文字列を意味する。

10

【0028】

図3の説明に戻る。文書データ群210bは、複数の文書データを含む。また、各文書データは、固有のIDが割り当てられ、各種の文字列を含む。

【0029】

単語インデックス210cは、文書データ群210bに含まれる各文書データの単語と、この単語の存在する文書データのIDとを対応付けた転置インデックスである。文字インデックス210dは、文書データ群210bに含まれる各文書データの文字と、この文字の存在する文書データのIDとを対応付けた転置インデックスである。

20

【0030】

インデクシング処理部220は、文書データ群210bから単語インデックス210cと文字インデックス210dを生成する処理部である。インデクシング処理部220は、単語区切方式により、文書データ群210bから単語インデックス210cを生成する。また、インデクシング処理部220は、文字区切方式により、文書データ群210bから文字インデックス210dを生成する。なお、単語区切方式による単語インデックス210cの生成は、周知の単語区切方式と同様である。文字区切方式による文字インデックス210dの生成は、周知の文字区切方式と同様である。

【0031】

入力受付部230は、利用者端末60から検索キーワードを受け付け、この検索キーワードを検索式解析処理部240に出力する。なお、入力受付部230は、検索装置200に接続された入力装置から、検索キーワードを取得してもよい。入力装置は、例えば、マウスやキーボードに対応する。

30

【0032】

検索式解析処理部240は、検索キーワードとパターンファイル210aとを比較して、単語インデックス210cを用いて文書データの検索を行うのか、文字インデックス210dを用いて文書データの検索を行うのかを判定する処理部である。以下において、単語インデックス210cを用いて文書データの検索を行うことを、単語区切方式の検索と表記する。文字インデックス210dを用いて文書データの検索を行うことを、文字区切方式の検索と表記する。

40

【0033】

まず、検索式解析処理部240は、検索キーワードに対して構文解析を実行する。例えば、検索キーワードを「1.1.1.1ANDsuppressor」とする。検索式解析処理部240が、検索キーワード「1.1.1.1ANDsuppressor」に対して構文解析を実行することで、この検索キーワードに含まれる条件文「AND」と、条件文を挟む文字列「1.1.1.1」、「suppressor」が抽出される。

【0034】

検索式解析処理部240は、検索キーワードから抽出した各文字列と、パターンファイル210aのパターンとをそれぞれ比較し、各文字列に対して、単語区切方式の検索を行

50

うのか、文字区切方式の検索を行うのかを判定する。

【0035】

具体的には、検索式解析処理部240は、パターンファイル210aのパターンのいずれかに文字列がマッチする場合には、文字区切方式の検索を行うと判定する。例えば、文字列「1.1.1.1」は、図4に示したパターンファイル210aのNo「2」のパターンとマッチする。このため、検索式解析処理部240は、文字列「1.1.1.1」に対して単語区切方式の検索を行うと判定する。

【0036】

また、文字列「suppressor」は、図4に示したパターンファイル210aのパターンとマッチしない。このため、検索式解析処理部240は、文字列「suppressor」に対して文字区切方式の検索を行うと判定する。

10

【0037】

検索式解析処理部240は、文字列と判定結果とを対応付けたデータをスコアリング処理部250に出力する。また、検索式解析処理部240は、検索キーワードに含まれる条件文も合わせてスコアリング処理部250に出力する。

【0038】

スコアリング処理部250は、検索式解析処理部240の文字列、文字列の判定結果、条件文を取得し、取得したデータに基づいて、検索キーワードに対応する文書データを検索する処理部である。ここでは一例として、文字列「1.1.1.1」に対応する判定結果が「単語区切方式の検索を行う」であり、文字列「suppressor」に対応する判定結果が「文字区切方式の検索を行う」であり、条件文が「AND」とする。

20

【0039】

この場合には、スコアリング処理部250は、文字列「1.1.1.1」と、文字インデックス210dとを比較して、文字列「1.1.1.1」に対応する文書データを特定し、特定した文書データを文書データ群210bから取得する。また、スコアリング処理部250は、文字列「suppressor」と、単語インデックス210cとを比較して、文字列「suppressor」に対応する文書データを特定し、特定した文書データを文書データ群210bから取得する。

【0040】

そして、スコアリング処理部250は、条件文が「AND」であるため、文字列「1.1.1.1」に対応する文書データと文字列「suppressor」に対応する文書データとを比較し、重複する文書データを検索結果出力部260に出力する。なお、条件文が「OR」の場合には、スコアリング処理部250は、条件文が「OR」であるため、文字列「1.1.1.1」に対応する文書データと文字列「suppressor」に対応する文書データとを検索結果出力部260に出力する。

30

【0041】

スコアリング処理部250は、文書データを検索した場合に、文書データに含まれる文字列の頻度に応じて、文書データにスコアを付与してもよい。

【0042】

検索結果出力部260は、スコアリング処理部250から受け付けた文書データを、利用者端末60に通知する。検索結果出力部260は、文書データのスコアに応じて、利用者端末60に表示させる文書データの順番を調整してもよい。また、検索結果出力部260は、検索装置200に接続された表示装置に文書データを出力してもよい。表示装置は、例えば、モニタや液晶ディスプレイに対応する。

40

【0043】

次に、本実施例2にかかる検索装置200の処理手順について説明する。図5は、本実施例2にかかる検索装置200の処理手順を示すフローチャートである。図5に示すように、検索装置200は、検索キーワードを取得し(ステップS101)、構文解析を実行する(ステップS102)。

【0044】

50

検索装置 200 は、パターンファイル 210 a から未選択のパターンを取得し（ステップ S103）、検索キーワードはパターンにマッチするか否かを判定する（ステップ S104）。検索装置 200 は、検索キーワードがパターンにマッチする場合には（ステップ S104, Yes）、単語区切方式の検索を行うと判定し（ステップ S105）、ステップ S108 に移行する。

【0045】

一方、検索装置 200 は、検索キーワードがパターンにマッチしない場合には（ステップ S104, No）、未選択のパターンが存在するか否かを判定する（ステップ S106）。検索装置 200 は、未選択のパターンが存在する場合には（ステップ S106, Yes）、ステップ S103 に移行する。

10

【0046】

一方、検索装置 200 は、未選択のパターンが存在しない場合には（ステップ S106, No）、文字区切方式の検索を行うと判定し（ステップ S107）、検索を実行する（ステップ S108）。

【0047】

上述してきたように、本実施例 2 にかかる検索装置 200 は、パターンファイル 210 a を用いて、単語区切方式の検索を行うのか、文字区切方式の検索を行うのかを判定している。このため、検索キーワードの特徴に合わせて最適なインデックスを選択することができるので、文書データの特性によらず、効率よく全文検索を実行することができる。

【0048】

また、本実施例 2 では、検索キーワードが検索式の場合に、この検索式を複数の部分キーワードに分割し、部分キーワード毎に単語区切方式の検索を行うのか、文字区切方式の検索を行うのかを判定している。このため、既存の技術を踏襲した検索式をそのまま利用して、全文検索を実行することができる。

20

【0049】

また、本実施例 2 のパターンファイル 210 a は、利用者単位の好みに合わせて容易にカスタマイズすることができる。

【0050】

ところで、図 3 に示した検索装置 200 の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、検索装置 200 の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。例えば、記憶部 210 を、着脱可能な外部装置または携帯端末等に搭載し、かかる外部装置または携帯端末等を検索装置 200 に有線または無線で接続するようにしてもよい。

30

【0051】

なお、検索装置 200 は、既知のパーソナルコンピュータ、ワークステーション、携帯電話、PHS 端末、移動体通信端末または PDA などの情報処理装置に、検索装置 200 の各機能を搭載することによって実現することもできる。

【0052】

図 6 は、実施例 2 にかかる検索装置を構成するコンピュータのハードウェア構成を示す図である。図 6 に示すように、このコンピュータ 300 は、各種演算処理を実行する CPU (Central Processing Unit) 301 と、ユーザからのデータの受け付ける入力装置 302 と、モニタ 303 を有する。また、コンピュータ 300 は、記憶媒体からプログラム等を読み取る媒体読み取り装置 304 と、ネットワークを介して他のコンピュータとの間でデータの授受を行うネットワークインターフェース装置 305 を有する。また、コンピュータ 300 は、各種情報を一時記憶する RAM (Random Access Memory) 306 と、ハードディスク装置 307 を有する。各装置 301 ~ 307 は、バス 308 に接続される。

40

【0053】

50

そして、ハードディスク装置 307 には、図 3 に示した検索式解析処理部 240、スコアリング処理部 250、インデクシング処理部 220 と同様の機能を有する検索プログラム 307a を記憶する。また、ハードディスク装置 307 は、図 3 に示した各種データ 210a ~ 210d にそれぞれ対応する各種データ 307b を記憶する。

【0054】

CPU 301 が検索プログラム 307a をハードディスク装置 307 から読み出して RAM 306 に展開することにより、検索プログラム 307a は、検索プロセス 306a として機能するようになる。また、CPU 301 は、各種データ 307b を RAM 306 に読み出す。検索プロセス 306a は、各種データ 306b を利用して、全文検索を実行する。

10

【0055】

なお、上記の検索プログラム 307a は、必ずしもハードディスク装置 307 に格納されている必要はなく、CD-ROM 等の記憶媒体に記憶されたプログラムを、コンピュータ 300 が読み出して実行するようにしてもよい。また、公衆回線、インターネット、LAN (Local Area Network)、WAN (Wide Area Network) 等にこのプログラムを記憶させておき、コンピュータ 300 がこれらからプログラムを読み出して実行するようにしてもよい。

【0056】

以上の各実施例を含む実施形態に関し、さらに以下の付記を開示する。

【0057】

20

(付記 1) 第 1 の区切方式に基づいて区切られ、文書データに関連付けられた第 1 のインデックスと、第 2 の区切方式に基づいて区切られ、文書データに関連付けられた第 2 のインデックスと、所定の文字の特徴を定義したパターンファイルとを記憶する記憶部と、

検索キーワードを受け付け、前記検索キーワードと前記パターンファイルとを基にして、前記第 1 のインデックスを用いて文書データの検索を行うのか、前記第 2 のインデックスを用いて文書データの検索を行うのかを判定する判定部と、

前記判定部の判定結果に基づいて、前記第 1 のインデックスまたは前記第 2 のインデックスを用いて文書データの検索を実行する検索部と

を備えたことを特徴とする検索装置。

【0058】

30

(付記 2) 前記検索キーワードを、複数の部分キーワードに分割するキーワード分割部を更に有し、前記判定部は、部分キーワード毎に前記第 1 のインデックスを用いて文書データの検索を行うのか、前記第 2 のインデックスを用いて文書データの検索を行うのかを判定することを特徴とする付記 1 に記載の検索装置。

【0059】

(付記 3) 前記第 1 の区切方式は、意味を持つ単語毎に文字列を区切る単語区切方式であり、前記判定部は、前記パターンファイルに定義された特徴が、前記検索キーワードと一致する場合には、前記第 1 のインデックスを用いて文書データの検索を行うと判定することを特徴とする付記 1 または 2 に記載の検索装置。

【0060】

40

(付記 4) 第 1 の区切方式に基づいて区切られ、文書データに関連付けられた第 1 のインデックスと、第 2 の区切方式に基づいて区切られ、文書データに関連付けられた第 2 のインデックスと、所定の文字の特徴を定義したパターンファイルとを記憶する記憶装置を有する検索装置が、

検索キーワードを受け付け、前記検索キーワードと前記パターンファイルとを基にして、前記第 1 のインデックスを用いて文書データの検索を行うのか、前記第 2 のインデックスを用いて文書データの検索を行うのかを判定する判定ステップと、

前記判定ステップの判定結果に基づいて、前記第 1 のインデックスまたは前記第 2 のインデックスを用いて文書データの検索を実行する検索ステップと

を含むことを特徴とする検索方法。

50

【 0 0 6 1 】

(付記5) 前記検索キーワードを、複数の部分キーワードに分割するキーワード分割ステップを更に含み、前記判定ステップでは、部分キーワード毎に前記第1のインデックスを用いて文書データの検索を行うのか、前記第2のインデックスを用いて文書データの検索を行うのかを判定することを特徴とする付記4に記載の検索方法。

【 0 0 6 2 】

(付記6) 前記第1の区切方式は、意味を持つ単語毎に文字列を区切る単語区切方式であり、前記判定ステップでは、前記パターンファイルに定義された特徴が、前記検索キーワードと一致する場合には、前記第1のインデックスを用いて文書データの検索を行うと判定することを特徴とする付記4または5に記載の検索方法。

10

【 0 0 6 3 】

(付記7) 第1の区切方式に基づいて区切られ、文書データに関連付けられた第1のインデックスと、第2の区切方式に基づいて区切られ、文書データに関連付けられた第2のインデックスと、所定の文字の特徴を定義したパターンファイルとを記憶する記憶装置を有するコンピュータに、

検索キーワードを受け付け、前記検索キーワードと前記パターンファイルとを基にして、前記第1のインデックスを用いて文書データの検索を行うのか、前記第2のインデックスを用いて文書データの検索を行うのかを判定する判定手順と、

前記判定ステップの判定結果に基づいて、前記第1のインデックスまたは前記第2のインデックスを用いて文書データの検索を実行する検索手順と

20

を実行させることを特徴とする検索プログラム。

【 0 0 6 4 】

(付記8) 前記検索キーワードを、複数の部分キーワードに分割するキーワード分割手順を更にコンピュータに実行させ、前記判定手順は、部分キーワード毎に前記第1のインデックスを用いて文書データの検索を行うのか、前記第2のインデックスを用いて文書データの検索を行うのかを判定することを特徴とする付記7に記載の検索プログラム。

【 0 0 6 5 】

(付記9) 前記第1の区切方式は、意味を持つ単語毎に文字列を区切る単語区切方式であり、前記判定手順は、前記パターンファイルに定義された特徴が、前記検索キーワードと一致する場合には、前記第1のインデックスを用いて文書データの検索を行うと判定することを特徴とする付記7または8に記載の検索プログラム。

30

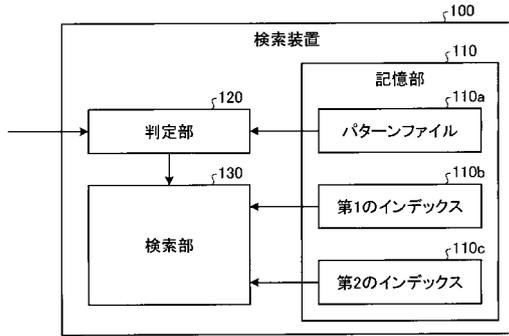
【符号の説明】

【 0 0 6 6 】

- 1 0 0 検索装置
- 1 1 0 a パターンファイル
- 1 1 0 b 第1のインデックス
- 1 1 0 c 第2のインデックス
- 1 2 0 判定部
- 1 3 0 検索部

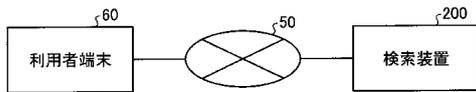
【図1】

本実施例1にかかる検索装置の構成を示す図



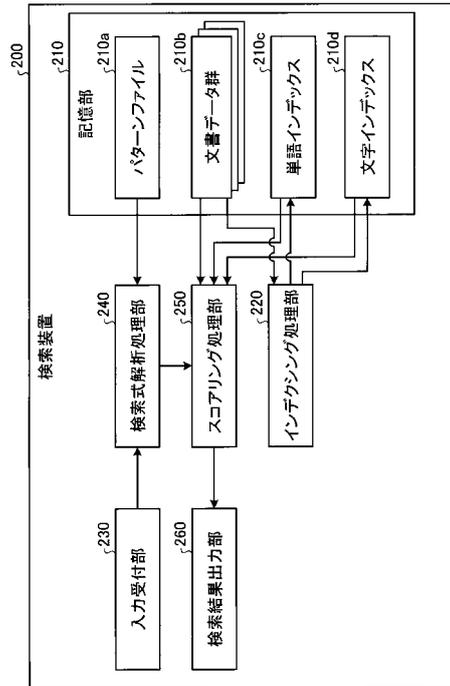
【図2】

本実施例2にかかるシステムを示す図



【図3】

本実施例2にかかる検索装置の構成を示す図



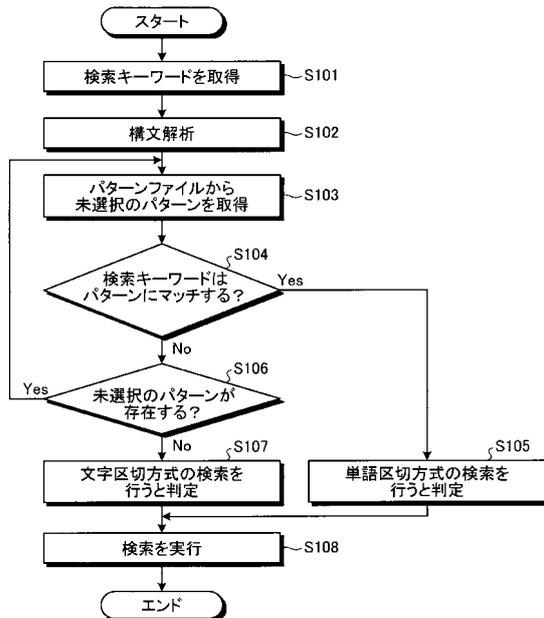
【図4】

パターンファイルのデータ構造を示す図

No	パターン
1	[0-9]{3}[0-9]*
2	[0-9]{1,3}[0-9]{1,3}[0-9]{1,3}[0-9]{1,3}
3	ENSG[0-9]+
4	NP_[0-9]+
5	IP[0-9]+
6	[a-z]{3,4}[a-zA-z0-9]+
...	...

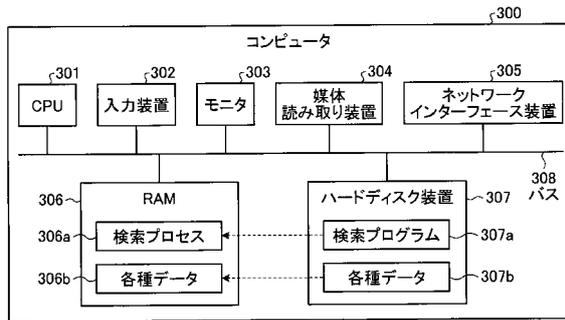
【図5】

本実施例2にかかる検索装置の処理手順を示すフローチャート



【図 6】

実施例にかかる検索装置を構成するコンピュータのハードウェア構成を示す図



フロントページの続き

(56)参考文献 特開平09 - 259132 (JP, A)
特開2008 - 077673 (JP, A)
特開2001 - 034623 (JP, A)
特開平10 - 307835 (JP, A)

(58)調査した分野(Int.Cl., DB名)
G06F 17/30
JSTPlus (JDreamIII)