

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3552318号

(P3552318)

(45) 発行日 平成16年8月11日(2004.8.11)

(24) 登録日 平成16年5月14日(2004.5.14)

(51) Int. Cl.⁷

G06F 17/30

F I

G06F 17/30 414B

G06F 17/30 170A

請求項の数 6 (全 12 頁)

(21) 出願番号	特願平7-2405	(73) 特許権者	000005108
(22) 出願日	平成7年1月11日(1995.1.11)		株式会社日立製作所
(65) 公開番号	特開平8-190571		東京都千代田区神田駿河台四丁目6番地
(43) 公開日	平成8年7月23日(1996.7.23)	(74) 代理人	100075096
審査請求日	平成13年3月15日(2001.3.15)		弁理士 作田 康夫
		(72) 発明者	島山 敦
			神奈川県川崎市麻生区王禅寺1099番地
			株式会社日立製作所 システム開発研究
			所内
		(72) 発明者	多田 勝己
			神奈川県川崎市麻生区王禅寺1099番地
			株式会社日立製作所 システム開発研究
			所内

最終頁に続く

(54) 【発明の名称】 文書検索方法およびシステム

(57) 【特許請求の範囲】

【請求項1】

文書検索方法において、
 記録装置に格納された文書ファイルから検索対象となる接続文字を含む文書ファイルの数を前記接続文字ごとに算出し、
 算出した前記接続文字を含む文書ファイルの数が、しきい値より大きい場合には、先頭のビットから文書ファイルの番号に対応したビット列を用いて文書識別子情報を作成し、
 算出した前記接続文字を含む文書ファイルの数が、しきい値より小さい場合には、前記接続文字が含まれる文書ファイルの番号のリストの形式で文書識別子情報を作成し、
 作成した前記文書識別子情報を記録装置に格納し、
 入力された検索語から接続文字を切り出し、
 切り出した前記接続文字それぞれについて前記文書識別子情報を読み出し、
 読み出した前記文書識別子情報の積集合を求めることにより、前記入力された検索語に含まれる接続文字を含む文書を抽出することを特徴とする文書検索方法。

【請求項2】

文書検索方法において、
 記録装置に格納された文書ファイルから検索対象となる接続文字を含む文書ファイルの数を前記接続文字の種類ごとに算出し、
 算出した前記接続文字を含む文書ファイルの数が、しきい値より大きい場合には、先頭のビットから文書ファイルの番号に対応したビット列を用いて文書識別子情報を作成し、

10

20

算出した前記接続文字を含む文書ファイルの数が、しきい値より小さい場合には、前記接続文字が含まれる文書ファイルの番号のリストの形式で文書識別子情報を作成し、作成した前記文書識別子情報を記録装置に格納し、前記文書識別子情報に基づいて文書検索を行うことを特徴とする文書検索方法。

【請求項3】

文書検索方法において、記録装置に格納された文書ファイルから検索対象となる接続文字の出現頻度を前記接続文字の種類ごとに算出し、算出した前記接続文字を含む文書ファイルの出現頻度が、しきい値より大きい場合には、先頭のビットから文書ファイルの番号に対応したビット列を用いて文書識別子情報を作成し、

10

算出した前記接続文字を含む文書ファイルの出現頻度が、しきい値より小さい場合には、前記接続文字が含まれる文書ファイルの番号のリストの形式で文書識別子情報を作成し、作成した前記文書識別子情報を記録装置に格納し、前記文書識別子情報に基づいて文書検索を行うことを特徴とする文書検索方法。

【請求項4】

文書検索方法において、前記コンピュータは、接続文字の第一文字目を配列の要素に含む文字テーブルを記録装置へ格納し、

前記接続文字の第二文字目と前記接続文字を含む文書ファイルに関する情報である文書識別子情報を格納したファイルへのポインタ情報とを含むファイルポインタテーブルを記録装置へ格納し、

20

前記文字テーブルに含まれる接続文字の第一文字目と当該接続文字の第二文字目とを対応づけるように前記文字テーブルから前記ファイルポインタテーブルへのポインタ情報を前記文字テーブルに格納し、

指定された検索語に含まれる接続文字について、前記検索語に含まれる接続文字の第一文字目に対応づけられた前記文字テーブルのポインタ情報を参照し、

参照した前記ポインタ情報に基づいて前記検索語に含まれる接続文字の第二文字目を格納するファイルポインタテーブルを参照して前記検索語に含まれる接続文字を含む文書識別子情報を読み出すことを特徴とする文書検索方法。

30

【請求項5】

文書検索方法において、記録装置に格納された文書ファイルを分割し、前記分割された文書ファイルごとに検索対象となる接続文字の文書識別子情報を作成し、前記分割された文書ファイルごとに作成した文書識別子情報をマージして記録装置に格納し、

格納した前記文書識別子情報に基づいて検索処理を行うことを特徴とする文書検索方法。

【請求項6】

文書検索システムにおいて、前記文書検索システムは、記録装置に格納された文書ファイルから検索対象となる接続文字を含む文書ファイルの数を前記接続文字の種類ごとに算出する手段と、

40

算出した前記接続文字を含む文書ファイルの数が、しきい値より大きい場合には、先頭のビットから文書ファイルの番号に対応したビット列を用いて文書識別子情報を作成する手段と、

算出した前記接続文字を含む文書ファイルの数が、しきい値より小さい場合には、前記接続文字が含まれる文書ファイルの番号をリストの形式で文書識別子情報を作成する手段と、

作成した前記文書識別子情報を記録装置に格納する手段と、

前記文書識別子情報に基づいて文書検索を行う手段とを含むことを特徴とする文書検索システム。

50

【発明の詳細な説明】

【0001】

【産業上の利用分野】

本発明は、文書データベースを、所定の文字列すなわち検索語を指定して文書の全文を対象として検索することにより、所望の文書を検索する文書検索方法に係わるものである。特に大量な文書を高速な検索を行う場合に好適な情報検索方法に関し、大規模文書データベースに適用されるものである。

【0002】

【従来技術】

先に、文書の登録の際にキーワード付けを行う必要のないフルテキストサーチ方式を特願平2-193015号(特開平3-174652号公報参照)で提案した。この方式は、文書を単語単位に圧縮した凝縮本文と、文書中の使用文字を一文字単位で登録した文字成分表を用いて、検索語に関連しない文書をふるい落とすことによってサーチ速度を等価的に高め、フルテキストサーチを実用レベルで高速に行うことを目的としたものである。また、この文字成分表を改良し更に高速なフルテキストサーチを実現する接続文字成分表方式を特願平3-342695号(特開平5-174064号公報参照)で提案した。この従来技術で用いる接続文字成分表は、テキストの中に含まれる所定の長さの接続する文字列を重複なく全て取り出し、これらを含む文書の識別子情報をビット列で記述するものである。しかし、全ての接続文字について識別子情報をビット列で記述すると、文字の組み合わせの個数分だけビット列が必要となり、接続文字成分表が膨大な容量になる。そこで、この従来技術では、ハッシュ関数を用いて1個のビット列に複数個の接続文字を割り当てておくようにして、容量を抑える工夫をしている。

【0003】

【発明が解決しようとする課題】

しかしながら、従来のハッシュ関数を用いて1個のビット列に複数個の接続文字を割り当てた場合には、同じビット列にまったく別の接続文字の文書識別子情報も重畳されることになる。従って、ある接続文字を指定して該当するビット列から文書識別子情報を取り出した場合、その情報からはまったく別の接続文字を含む文書が得られる可能性がある。つまり、ハッシュ関数を用いた接続文字成分表による検索結果には検索ノイズが含まれることになる。このことは、大量の文書を登録する大規模な文書検索システムでは、検索語に関連しない不要な文書のふるい落とし、すなわち絞り込みが適切に行われないう可能性があることを意味し、その場合には検索性能の低下につながる。

【0004】

ハッシュ関数を用いずに、全ての接続文字についてそれぞれ1個のビット列を対応させることも考えられるが、その場合にはビット列のデータ量が膨大なものとなり、実用的ではない。具体的に説明すると、日本語で使用する文字コードは、現在約8,000種類あるので、2文字の組み合わせとしての接続文字の種類は、 $8,000 \times 8,000 = 6,400$ 万種類となる。登録する文書数を100万件とした場合、この6,400万種類のそれぞれの接続文字に100万bitの文書識別子情報を対応させるので、 $6,400$ 万種類 \times 100万bit = 8TByteもの容量が必要になる。この文字成分表の大きさに対し、文書本体の大きさを20KB/件としても、100万件で、 $20KB \times 100$ 万件 = 20GByteであり、圧倒的に文字成分表の容量のほうが大きくなってしまふ。

【0005】

すなわち、本発明の解決しようとする課題は、大規模な情報検索システムにおいても検索ノイズの少ない接続文字成分表を、実用的な容量で実現することにある。

【0006】

【課題を解決するための手段】

本発明は、以下の構成を採ることにより上述の課題を解決する。

【0007】

文書のテキストデータにおける複数の文字の共起関係を記述した接続文字を接続文字ファ

10

20

30

40

50

イルに重複なく格納する接続文字格納ステップと、前記接続文字ファイルに格納された接続文字を参照して、指定した条件式中の検索語に含まれる接続文字を含む文書を検索結果の候補とする文書検索方法において、接続文字格納ステップとして、テキストデータ中に現れる接続文字成分の種類および各接続文字成分の出現する文書数を算出し、算出された文書数が所定のしきい値より大きい場合は該当文書の文書番号に対応する位置を“1”とするビット列として登録し、しきい値より小さい場合には該当文書の文書番号をバイナリデータとして格納することを特徴とする。

【0008】

より詳細に言うると以下の(1)～(6)の各ステップに分けることができる。

【0009】

(1) テキストデータ分割ステップ
 (2) 文書識別子情報作成ステップ
 (3) 文書識別子情報マージステップ
 (4) 検索語分割ステップ
 (5) 文書識別子情報探索ステップ
 (6) 文書識別子情報ANDステップ
 (1)から(3)は文字成分表の登録のための処理であり、(4)から(6)はこれを用いた検索のための処理である。これより、各ステップの処理内容を説明する。

【0010】

(1) テキストデータ分割ステップ
 文字成分表への登録の際、文字の組合せの個数および各組合せに対応する文書識別子の記憶容量を抑えるために一回に処理する文書数を適切な数に分割する。分割する文書数は、予め設定してもよいし、登録に使用する計算機のメモリ容量から算出してもよい。

【0011】

(2) 文書識別子情報作成ステップ
 (1)で分割した文書群のそれぞれについて別個に文書識別子情報を作成していく。具体的には、文書中に実際に現われた文字の組合せとその文字の組合せが現われた文書識別子の情報を対にして格納する。

【0012】

(3) 文書識別子情報マージステップ
 (2)で作成した文書識別子情報を(1)で分割した文書群の数分マージして、登録文書全体の文字成分表を作成する。

【0013】

(4) 検索語分割ステップ
 与えられた検索語を登録時と同じ方法で文字の組合せに分割する。

【0014】

(5) 文書識別子情報探索ステップ
 (4)で分割した文字のそれぞれについて、文書識別子情報を探索する。

【0015】

(6) 文書識別子情報ANDステップ
 (5)で得られた文書識別子情報のそれぞれについて、AND処理を行うことにより、与えられた検索語の全ての接続文字を含む文書を文字成分表サーチ結果として出力する。

【0016】

【作用】

以下、これらのステップからなる本発明の文書検索方法の原理を説明した上で、その作用を説明する。

【0017】

まず、本発明で用いる文字成分表の構成について説明する。本発明では、接続文字に対応する文書識別子情報を管理するのに、文字テーブル、ファイルポインタテーブルを用いる。図2は文字テーブルおよびファイルポインタテーブルの概要を示す図である。

10

20

30

40

50

【 0 0 1 8 】

たとえば、“構成”という文字列を含む文書を検索する場合には、まず文字テーブルについて“構”の文字に対応するレコードを参照してファイルポインタテーブルへのポインタ情報580を得る。次に、ファイルポインタテーブルの先頭から580バイト目からの各レコードを参照して、第二文字目が“成”のレコードを探索する。ファイルポインタテーブルには、各接続文字の第一文字目ごとに、先頭に第二文字目が0のレコードを格納しておく。第二文字目が0のレコードには、第一文字目の一文字を含んでいる全ての文書の文書識別子情報へのポインタを格納しておく。すなわち、第二文字目が0のレコードは、第一文字だけからなる単一文字に対応する文書識別子情報をアクセスするためのファイル識別子（以後ファイルIDとも呼ぶ）とファイル内バイト位置（以後オフセットとも呼ぶ）を格納する。したがって、各接続文字ごとに第二文字目が0のレコードが必ず存在するため、例えば、“構成”の接続文字を探索する場合は、“構”に対応するファイルポインタテーブルの先頭から580バイト目のレコードから探索を開始し、再び第二文字目が0になるまで探索を続け、もし“成”の文字が見つからない場合は、該当する接続文字がないと判断できる。図2の例では、“成”のレコードが存在するため、ここからファイルIDが1、オフセットが1034という文書識別子情報へアクセスするための情報を得ることができる。

10

【 0 0 1 9 】

文書識別子情報は、図3のように複数のファイルに分割格納する。ファイルポインタテーブルのファイルID情報により、どのファイルに文書識別子情報が格納されているかを特定する。なおかつ特定のファイルIDは、文書識別子情報をビット列で持つとあらかじめ決めておく。図3の例では、ファイル1が文書識別子情報をビット列で持つファイルとしている。図2の例では、接続文字“構成”に関する文書識別子情報へのアクセス情報として、ファイルIDが1、オフセットが1,034が得られる。したがって、ファイル1内の1,034バイト目からのビット列“0111010101...”が文書識別子情報として得られることになる。このビット列は、先頭ビットから文書番号に対応して、“1”が接続文字“構成”を含む文書を示すことになる。すなわち、この例では、“構成”を含む文書の文書番号は、1、2、3、5、7、9...となる。図3の他のファイル（ファイル2及びファイル3）は文書識別子情報をIDリストの形式で格納したものである。各IDリストの先頭は格納してある文書番号の個数を示している。例えば、接続文字“構造”の場合、図2の例では、ファイルIDが2、オフセットが340であるので、ファイル2の先頭から340バイト目を参照することによって、接続文字“構造”を含む文書数が56個あり、文書番号が562、1038、...であることがわかる。

20

30

【 0 0 2 0 】

このように、ファイルポインタテーブルには、データベース中に存在する接続文字のみを登録するので、データベース中に存在しない文字の組み合わせは全て排除できるという利点がある。したがって、文字テーブルやファイルポインタテーブルで実現している接続文字の管理情報を格納するファイル量やメモリ量を大幅に削減することができる。また、文書識別子情報をビット列あるいはIDリストの形式で格納し、多くの文書を格納する場合はビット列で、少ない文書を格納する場合はIDリストの形式で管理することによりファイル容量を大幅に削減することができる。具体的に説明すると、ビットリストの形式で文書識別子情報を格納するには、常にデータベースに登録した全件分のビット数が必要になるが、IDリストの形式で文書識別子情報を格納する場合には、文書識別子を表わすビット数×登録文書数ですむことになる。例えば、データベースの全登録件数が100万件で、一個の文書識別子情報を表わすのに32ビットを割り当てるとすると以下の格納領域が必要となる。接続文字“構造”を含む文書を10件登録する場合に、ビット列ならば、100万bit = 125KBの格納領域が必要となるが、IDリスト形式ならば、32bit × 10件 = 40Bの格納領域ですむことになる。一方、例えば、接続文字“構成”を含む文書が100万件中で90万件ある場合には、ビット列ならば、100万bit = 125KBの格納領域にすむのに対し、IDリスト形式の場合、32bit × 90万件 = 3 .

40

50

6MBの領域が必要となる。したがって、この100万件を、文書識別子32ビットで格納する場合には、 $100万\text{bit} \div 32\text{bit} = 31,250$ 件を境として、これよりも登録件数が多い場合はビット列形式で、少ない場合はIDリスト形式で文書識別子情報を格納するのが、最も格納領域を有効に使用方法である。

【0021】

次に、このような文字成分表の登録の方法について、原理を説明する。文字テーブルとファイルポインタテーブルを用い、データベース中に用いられる接続文字のみを文字成分表に登録することにより、ファイル容量を実用容量に抑えることができることは既に説明した。

【0022】

したがって、登録時に全ての接続文字成分について管理をしようとする、メモリ容量が足りなくなり、文字成分表を作ることが不可能となる。磁気ディスクをワークにして情報を一旦退避する方法もあるが、アクセス速度が遅いので登録処理に極めて時間が掛かることになる。そこで、図4のように登録するテキストデータを分割して、分割したテキストデータ毎に文字成分表を作成し、最後にこれらをマージして全テキストデータの文字成分表を作成する。図4では、全部で2万4千件のテキストデータを8千件毎に分割して文字成分表を作成する例を示している。“構成”という接続文字について、最初の8千件のテキストデータでは、文書番号50、145、290...が文書識別子情報として蓄えられる。同様に、次の8千件、その次の8千件についても各分割したテキストデータ毎に文字成分表を作成する。最後に、それぞれで得られた文書識別子情報をマージして、本図の例では、“構成”の接続文字に対する文書識別子情報として、50、145、290、8096、12365、17851、22989...という情報を作成する。

【0023】

検索の際には、入力された検索語を接続文字に分割し、それぞれの接続文字に対応する文書識別子情報を読み出してきて、それらの情報の積集合を取り、これを文字成分表の検索結果とする。すなわち、“建造物”という検索語については、“建造”と“造物”の2種類の接続文字について、それぞれ文字成分表の文書識別子情報を読み出してそれらの積を演算する。例えば、接続文字“建造”に対応する文書識別子情報が562、1038、2458...で、接続文字“造物”に対応する文書識別子情報が261、562、2458...の場合は、検索語“建造物”の文字成分表サーチ結果は文書番号で562、2458...となる。

【0024】

このように、各接続文字に対する文書識別子情報はノイズのない情報であるため、これらの文書識別子情報を論理式演算(AND)して得られる文字成分表サーチ結果も、従来のハッシングを行う文字成分表のサーチ結果に比べ、ハッシングに起因するノイズが除去されることになり、検索精度が大幅に向上できることになる。

【0025】

【実施例】

以下、本発明の実施例について図を用いて詳細に説明する。

【0026】

図1は、本実施例の構成を示す図である。本実施例は、登録検索用の端末101、102、...110、ネットワーク200、文書サーバ1000からなる。文書サーバ1000には、LANアダプタ1010、CPU1020、ワークメモリ1030、文字テーブル1100とファイルポインタテーブル1200を格納するメモリ、テキストデータ分割プログラム1310、文書識別子情報作成プログラム1320、文書識別子情報マージプログラム1330、検索語分割プログラム1340、文書識別子情報探索プログラム1350、文書識別子情報ANDプログラム1360を格納するメモリ、文字成分表を分割して格納するファイル1401、1402、...、テキストデータ1410からなる。

【0027】

データの登録時には、テキストデータ分割プログラム1310で登録する文書データを一

10

20

30

40

50

定の件数に分割し、分割したそれぞれのテキストデータについて文書識別子情報作成プログラム 1320 で文書識別子情報を作成して、最後に分割して作成したそれぞれの文書識別子情報を文書識別子情報マージプログラム 1330 でマージして文字テーブル 1100、ファイルポインタテーブル 1200、文字成分表 1401、1402、1403 を作成する。

【0028】

また、データの検索時には、各端末から与えられた検索語を検索語分割プログラム 1340 によって文書識別子情報を作成したときと同じアルゴリズムで接続文字に分割し、それぞれの接続文字について文書識別子情報探索プログラム 1350 で該当する文書識別子情報を文字成分表 1401、1402、1403 から取り出す。そして、検索語を構成する
10
全ての接続文字に対応する文書識別子情報を文書識別子情報 AND プログラム 1360 によって AND することで検索語を含む文書を文字成分表のサーチ結果とする。

【0029】

まず、データの登録処理に従い、文字成分表の作成手順を説明し、次に検索処理に従って文字成分表による候補文書の抽出過程を説明する。作用の項でも説明したように、大量の文書について文字成分表を一度に登録するには、大量のメモリを使用しなければならないので、本実施例では 8,000 件ごとに小さな文字成分表を作成し、最後に一つの文字成分表に統合する処理を行う。図 5 に、この文書識別子情報作成処理の手順を示す。まず、8,000 件のそれぞれの文書について (5010) 接続文字の抽出 (5020) を行い、切り出した接続文字についてその出現頻度情報を計数 (5030) する。そして、算出
20
した出現頻度にしたがって文書識別子情報を格納するメモリアreaをワークメモリ上に確保し、それぞれの接続文字の出現頻度情報が所定のしきい値より大きい場合にはビット列で各接続文字が出現する文書番号を文書識別子情報として登録 (5040) していく。8,000 件の全ての文書について文書識別子情報を登録し終わったらファイルに文字テーブル、ファイルポインタテーブル、文書識別子情報を格納 (5050) しメモリ領域を解放する。8,000 件単位にこのように小さな分割文字成分表を作成し、最後に各分割文字成分表をマージ (5060) してデータベース全体の文字成分表を作成する。

【0030】

この分割文字成分表のマージ処理 (5060) は、図 6 に示すとおり各分割文字成分表の文字テーブルとファイルポインタテーブルを参照し、それぞれの接続文字に対応する文書
30
識別子情報を統合する形で進めていく。図 6 は二個の分割文字成分表を一個の文字成分表に統合する例を示している。具体的な処理の手順を図 7 に示す。まず、それぞれの分割文字成分表の文字テーブルを参照 (7010) し、統合した文字テーブルを作成 (7030) する。この時文字テーブルの各レコードについて (7020)、どちらか一方にしか登録されていないレコードについては、登録されている側に記録されているファイルポインタテーブルの各文字について (7040) 内容を統合したファイルポインタテーブルに登録する (7050) とともに、ファイルポインタテーブルで管理されている文書識別子情報をマージ前の小さな文字成分表からマージ後の文字成分表へコピー (7060) していく。また、双方の文字テーブルに同じ文字が存在する場合には、記録されているファイル
40
ポインタテーブルの各文字について (7070)、ファイルポインタテーブルに記載された第二文字目を比較しながら統合したファイルポインタテーブルを作成 (7080) していく。すなわち、ファイルポインタテーブルの第二文字目が一致しない場合には、該当する文書識別子情報をコピー (7090) し、一致する場合には双方の文書識別子情報をマージ (7100) して格納する。

【0031】

この文書識別子情報のマージ及びコピーの際には、マージ後の登録件数から所定の件数よりも多い場合にはビット列に、少ない場合には ID リストの形式にして格納する。

【0032】

以上のマージ処理アルゴリズムを図 6 を用いて具体的に説明する。“構”の文字は文字テ
50
ーブル 1 および文字テーブル 2 のどちらにも存在する。したがって、“構”の文字に対応

するファイルポインタテーブル1の内容とファイルポインタテーブル2の内容を統合ファイルポインタテーブルに登録していく。ファイルポインタテーブルにおける該当レコードの先頭の第二文字目が0のレコードは、“構”の一文字を含む文書の識別子情報をアクセスするための情報を格納している。この第二文字目が0のレコードはファイルポインタテーブル1とファイルポインタテーブル2の両方に存在するので、双方のファイルIDとオフセットで与えられる文書識別子情報をマージして統合文字成分表に登録する。“構”に対応するファイルポインタテーブルの第2レコード“成”についても同様である。第3レコードについてはファイルポインタテーブル1が“造”であるのに対して、ファイルポインタテーブル2では“築”と異なっている。したがって、それぞれの文書識別子情報をマージ前の小さな文字成分表から統合文字成分表へコピーする。

10

【0033】

検索処理は、図8に示す手順で行う。まず、検索語から接続文字を切り出す(8010)。次に、切り出したそれぞれの接続文字について(8020)、文字テーブルを探索する(8030)。そして、該当するファイルポインタテーブルの各レコードについて、第二文字目の探索を行い(8040)該当するファイルIDとオフセットを得る。こうして、得られた文書識別子情報を格納したファイルとそのオフセット値より、該当する接続文字に対応するIDリストまたはビット列を読み出し、IDリストの場合にはこれをビット列に変換することにより文書識別子情報を取得する(8050)。この文書識別子情報の取得の過程で該当する接続文字が文字成分表に登録されていない場合(8060、8070)には、すなわち検索語を構成する接続文字のうちどれか一つでも文字成分表に登録されていなければ、検索語を含む該当文書がないということを意味することになるため検索結果として0件という結果を、文書識別子情報探索プログラム1350がLANアダプタ1010を介して検索端末に返す。

20

【0034】

検索語を構成する全ての接続文字について該当する文書識別子情報が得られた場合は、得られたそれぞれの文書識別子情報の積集合をとることによって、指定された検索語中の全ての接続文字を含む文書のみを抽出することができる。

【0035】

このようにして得られた文字成分表の検索結果は、検索ノイズが非常に少ないので、文字成分表のサーチ結果を表示しても十分実用できる。もちろん、文字成分表のサーチ結果をもとに、文書本文を検索し実際に検索語を含む文書のみ絞り込むかあるいは、複数の検索語間の位置的関係を満たす文書を探すことも可能である。また、文字成分表の検索結果を一度検索端末に表示し、ユーザの指定により本文の探索を行うかどうかを決定してもよい。

30

【0036】

以上、本実施例によれば、データベース中に存在する接続文字のみを登録するので、データベース中に存在しない文字の組み合わせは全て排除できるという利点がある。また、文書識別子情報をビット列とIDリストの形式で格納し、多くの文書識別子情報を格納する場合はビット列で、少ない文書識別子情報を格納する場合はIDリストの形式で格納することでファイル容量を大幅に削減することができる。

40

【0037】

さらに、各接続文字に対して文書識別子情報は必ずその接続文字を含むノイズのない情報であるから、これらの文書識別子情報をANDして得られる文字成分表サーチ結果も、検索精度を大幅に向上することができる。

【0038】

また、本発明によれば2文字以上の接続文字についても登録することにより、さらに文字成分表サーチの検索ノイズを少なくすることも可能である。

【0039】**【発明の効果】**

本発明によれば、文書識別子情報をビット列とIDリストのどちらかの形式で選択的に格

50

納することにし、多くの文書識別子情報を格納する場合はビット列で、少ない文書識別子情報を格納する場合はIDリストの形式で格納することでファイル容量を大幅に削減することができる。

【0040】

また、各接続文字に対して文書識別子情報は必ずその接続文字を含むノイズのない情報であり、これらの文書識別子情報を検索語の接続文字の個数分ANDするので、文字成分表サーチの検索精度を大幅に向上することができる。これにより、検索語間の位置的な条件などを検索する場合にも、より本文情報の検索範囲を狭めることができるという利点がある。

【0041】

さらに、文字テーブル及びファイルポインタテーブルを用いることにより、データベース中に存在する接続文字のみを登録するので、データベース中に存在しない文字の組み合わせは全て排除できるので、接続文字を管理するために必要なメモリ量を少なくできるという利点がある。

【0042】

さらにまた、文字成分表の登録の際に、登録文書を分割して小さな分割文字成分表を作成し、後でこれらの分割文字成分表をマージして目的の文字成分表を作成することにより、少ないメモリ容量でも効率的に大きなデータベースの文字成分表を作成することができる。

。

【図面の簡単な説明】

【図1】本発明の第一の実施例の構成図である。

【図2】文字成分表のテーブル構成図である。

【図3】文書識別子情報格納ファイルの概要を示す図である。

【図4】文字成分表登録処理の概要を示す図である。

【図5】登録処理の流れを示すPAD図である。

【図6】分割文字成分表の統合処理を示す概念図である。

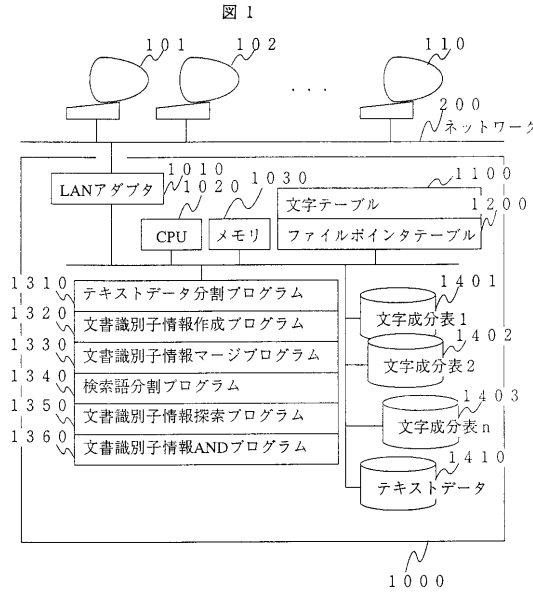
【図7】統合処理の流れを示すPAD図である。

【図8】検索処理の流れを示すPAD図である。

10

20

【 図 1 】



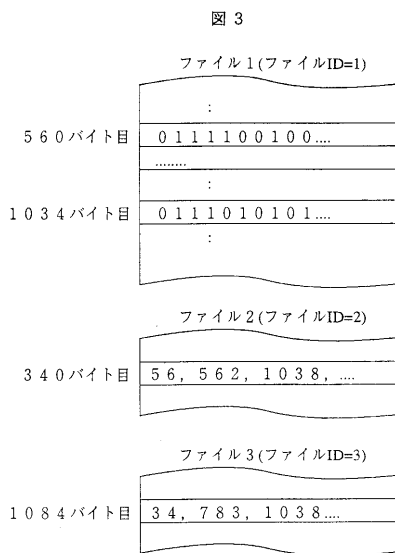
【 図 2 】

図 2

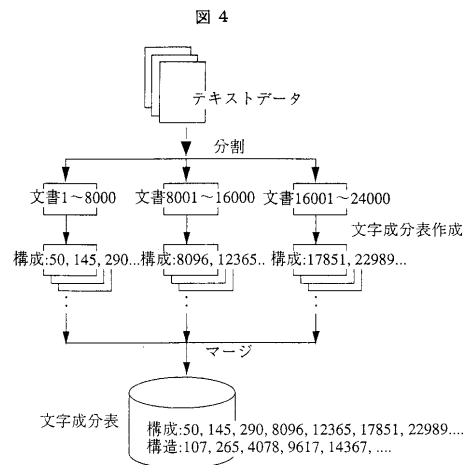
文字テーブル		ファイルポインタテーブル		
ポインタ		第二文字目	ファイルID	オフセット
:	:	:	:	:
構	580	0	1	560
:	:	成	1	1034
:	:	造	2	340
:	:	:	:	:
内	870	0	1	2460
:	:	部	3	1084
:	:	:	:	:

1100 1200

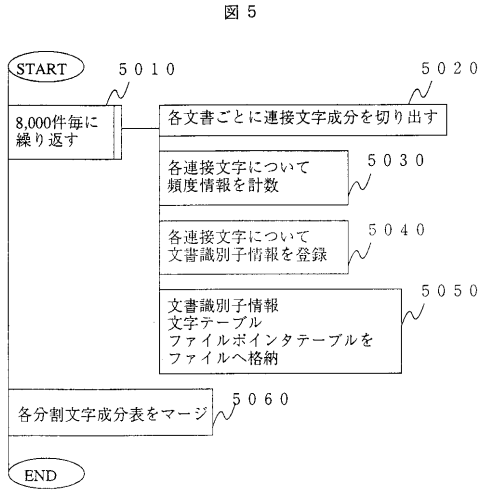
【 図 3 】



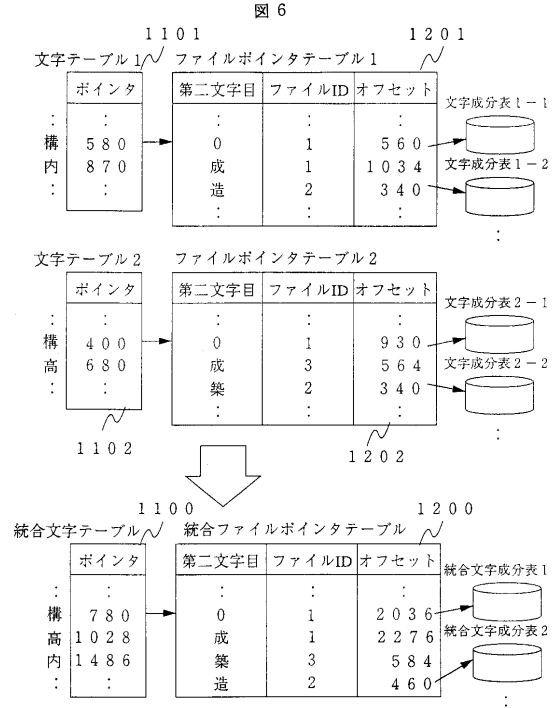
【 図 4 】



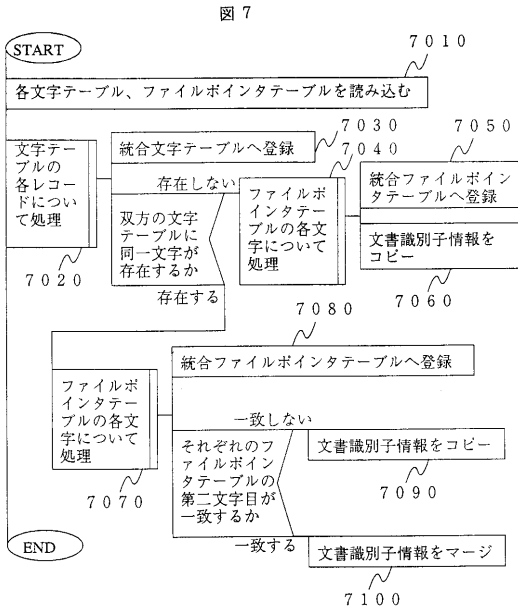
【 図 5 】



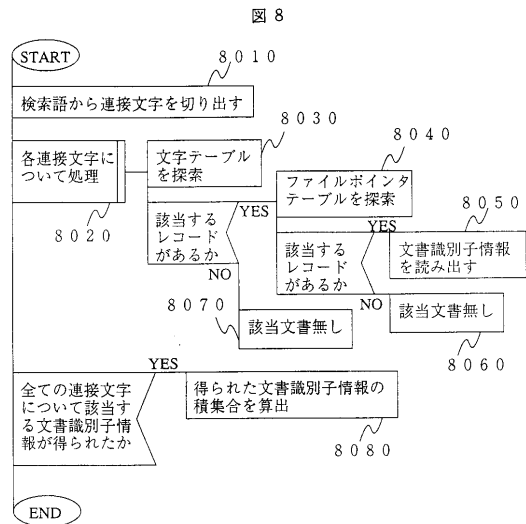
【 図 6 】



【 図 7 】



【 図 8 】



フロントページの続き

(72)発明者 加藤 寛次

神奈川県川崎市麻生区王禅寺1099番地 株式会社日立製作所 システム開発研究所内

(72)発明者 浅川 悟志

神奈川県横浜市戸塚区戸塚町5030番地 株式会社日立製作所 ソフトウェア開発本部内

審査官 田川 泰宏

(56)参考文献 特開平04-274557(JP,A)

特開平05-174064(JP,A)

特開平07-319920(JP,A)

特開平08-030633(JP,A)

岩崎雅二郎,小川泰嗣,文字成分表による文字列検索の実現と評価,情報処理学会研究報告(93-DBS-92),1993年3月22日,Vol.93,No.29,p.1-10

小川泰嗣,岩崎雅二郎,林大川,全文検索のための文字成分表方式の改良,情報処理学会研究報告(94-DBS-99),1994年7月22日,Vol.94,No.62,p.261-264

畠山敦,ソフトウェアによるテキストサーチマシンの実現,情報処理学会研究報告(92-FI-25),1992年5月12日,Vol.92,No.32,p.19-25

(58)調査した分野(Int.Cl.⁷,DB名)

G06F 17/30

JICSTファイル(JOIS)