

7H-7 キーワードの登録回数に基づくグループ化による一検索方式

菊池 忠一(1) 藤村 浩(2) 藍沢 實(1)
 (1) (株) テレマティーク国際研究所 (2)電気通信大学

1 はじめに

電子メディアの普及と通信の高度化に伴い、電子図書館の構築が進められるようになってきた。図書には、著者/発行者/件名等の数種のキーワードがあり、検索時にはこれらの中から任意のキーワードを使用されるのが普通である。また、一般的な数万冊から約100万冊蔵書を有する図書館の場合、現状、検索に時間がかかり、検索を煩わしいものになっている。

現在図書検索に用いられている検索手法としては次のものがある。①逐次サーチ法は検索ファイルのデータ量に比例した検索時間を要し、大量情報向きではない。②転置ファイル法は、検索入力数に比例した検索時間を要し、複数検索入力向きではない。これらの改良方式として、特性ファイル法のフォルスドロップを改良した重ね合わせ符号法がある。これは、システムの保守・管理が転置ファイルより容易で、逐次サーチ法より高速検索できる[1]が、大量情報検索には必ずしも十分とは言えない。

本報告では、大量情報における高速検索方式、特に、複数の検索入力における論理積条件下での高速検索方式を提案している。4万冊の書誌情報を対象とする検索実験で、平均プログラム実行時間2.3msの検索速度を確認した。

2 アルゴリズムの特徴

本方式は、複数のキーワードを使用し、大量情報の中から高速検索できるもので、以下の特徴を持つ。

- (1) 登録使用頻度の少ないキーワードから検索する方が高速性が得られる特徴を活かした検索ファイル構造である。
 - ① 検索ファイルには、キーワード登録使用頻度の低頻度から昇順に、全キーワードに対応した検索レコードがわりあてられる。
 - ② 各登録情報の有するキーワードについて、検索レコードに対応するキーワードの頻度より高頻度のキーワードをすべて格納する。(図1)
- (2) 検索ファイルは(1)の①により、各検索レコード容量が均等化する方向に作成される。

- (3) 複数のキーワードによる検索の場合、低頻度のキーワードに対応する検索レコードだけで、検索可能である。
- (4) キーワード種ごとにマーク付けし、すべて4バイトでコード化することにより、キーワード種別処理が不要である。

登録情報	検索ファイル	
j	kw1	
kw5		
Λ	kw5	j(kw5, kw10, kw25, kw100)
kw10		
Λ	kw10	j(kw10, kw25, kw100)
kw25		
Λ	kw25	j(kw25, kw100)
kw100		
	kw100	j(kw100)

図1 キーワード登録例

3 登録

3.1 キーワード表の作成

著者名と発行者名にそれぞれ専用マークを付加し、件名も含めた3種のキーワードを登録順に、1から昇順の数値を付与してコード化する。コード化と同時に登録回数を計数し、キーワード表(図4)を作成する。例えば図4では、ISDNはキーワードとして5番目に登録され、既にISDNをキーワードに持つ書誌情報が100件あったことを示す。

3.2 検索ファイルの作成

検索ファイル(図2)に、キーワード登録回数の最も少ないものから昇順に、すべてのキーワードに対応した検索レコード(図3)を割り当て、登録情報の有するキーワードの中から、検索レコードに対応するキーワードの登録回数より多いキーワードをすべて格納する。例えば、ISDN、通信、OSIを有する書誌番号10の場合、キーワード表から登録回数順はOSI < ISDN < 通信であるから、OSIのコード150に対応する検索レコードにISDN、通信、OSIを格納し(図3)、ISDNの

A Group Retrieval Technique Based on Registration Frequency of Keyword

Chuichi Kikuchi(1) Hiroshi Hujinura(2) Minoru Aizawa(1)

(1)Telematique Intl. Res. Lab. (2)University of Electro-Communications

コード5に対応する検索レコードにはISDNと通を格納し、通信のコード20に対応する検索コードには通信を格納する。

4 検索

キーワード表を使用して、入力されたキーワードの内、登録回数の最も少ないコードをレコード番号とする検索レコードから、入力されたすべてのキーワード(コード)が共通に持つ書誌番号を取り出す。

例えば、図3と図4において、検索入力が入力ISDN、通信、OSIの場合、登録回数が最も少ないOSIに該当するコード150をレコード番号として、検索レコードを取り出す。この検索レコードの中から、ISDN、通信、OSIのコード5、20、150に共通する書誌番号として10を取り出す。

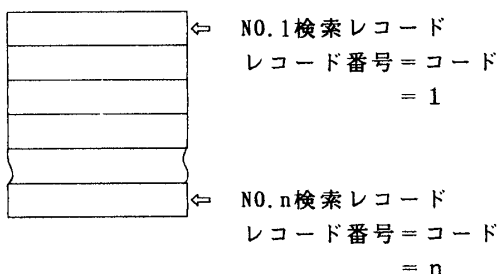


図2 検索ファイル

レコード番号	150
コード 5	書誌番号 10
5	7 8
20	3
20	1 0
20	2 2 2
150	1 0

図3 検索レコード

キーワード	コード	登録回数
ISDN	5	100
ソフト	17	35
通信	20	350
OSI	150	50

図4 キーワード表

5 性能試験

表記方式を評価するために、PFU社製A-50ミニコン(OS:UNIX)上に、40036冊の書誌情報とキーワードを登録した実験システムを構築し、登録時間、検索ファイル容量、書誌情報が有する全キーワード検索の時間を調べた。

使用した書誌情報(表1参照)は、キーワード(著者名、発行者名、件名)が1個の書誌情報が64件、2個の書誌情報が1218件、...、12個の書誌情報が33件である。キーワードは合計48008種あり、1書誌情報平均4.8個である。

試験結果は、検索ファイルは約4.4MB、キーワード表は約1.7MB、検索ファイル・ポインタ表は約0.5MBであった。登録時間は約5.6時

間、検索時間は表1に示すとおりであった。なお、検索時間は、検索ファイルから該当する全書誌番号を取り出しまでの時間で、ディスク装置アクセス時間を除くプログラム処理時間である。

検索入力 キーワード数	検索時間			検索回数 書誌情報数
	平均 ms	最小 ms	最大 ms	
1	1.5	1.1	8.8	64
2	1.6	1.3	5.9	1218
3	1.8	1.4	6.8	6324
4	2.0	1.6	11.8	12401
5	2.3	1.8	12.9	9422
6	2.7	1.9	12.9	5109
7	3.0	2.1	11.4	2865
8	3.3	2.3	10.2	1500
9	3.4	2.5	7.7	656
10	3.3	2.7	5.7	304
11	3.4	2.9	5.1	140
12	3.5	3.0	4.9	33
全体	2.3	1.1	12.9	40036

表1 検索時間

6 おわりに

約4万冊書誌情報を対象とする性能試験で、平均検索時間2.3msを確認できた。これはプログラム実行時間のみで、実際にはこれにディスク装置アクセス時間(約30ms)を加えたものになる。

表1の検索結果から平均検索時間と最小検索時間は、検索入力されたキーワード数nの一次関数で表わせることがわかる。

$$\text{平均検索時間 } t = 0.26(n-1) + 1.3 \text{ ms } n \leq 8$$

$$\text{最小検索時間 } t = 0.18(n-1) + 1.08 \text{ ms } n \leq 12$$

係数は、該当検索レコード内で検索入力キーワードに共通する書誌番号を抽出する論理積にかかる時間である。

平均検索時間に着目すると、n ≥ 8まではほぼ比例するが、n ≥ 9では飽和してくる。これは、キーワード8個前後で検索結果が出ていることと、書誌情報が急減していることによると考えられる。

本方式は、複数の検索入力に対し、ディスク装置から1個の検索レコードを取り出すだけで検索可能なことから、複数のキーワードを持つ大量情報の高速検索に有効な方式であると考えられる。

今回は、約4万冊を対象に性能試験を行い検索の高速性を確認したが、今後、本方式の特徴を明確にするために大量情報の検索を試みたい。

参考文献

[1] 有川、篠原他：「重ね合わせ符号を用いた文献検索システムについて」：情報研報データベースシステム54-2