

2N-5

キーワードのコード化による検索方式

菊池 忠一 飯島 豊  
(株) テレマティーク国際研究所

1. はじめに

最近、電子メディアの普及が著しいが、情報検索時に入力するキーワード(以下KWと略す)がわからないとか、KWの一部しかわからないという場合が多く、検索を煩わしいものになっている。

検索時に書名や見出し語の一部、時には単語の一部しか記憶していない場合があり、この場合、書名や見出し語の一部でも検索可能な部分一致や中間一致検索が有効になる(部分一致=前方一致U中間一致U後方一致)。しかし、現行の方式では転置ファイルが大きくなるために、部分一致や中間一致検索は利用できないことが多い。この解決の一手法として、インデックスを用いない国語辞書全文検索方法[1]が提案されているが、この手法は検索ファイル内の全キーワードを照合する必要がある。本報告では、KWをコード化して取り扱い、該当文字ブロックを照合することにより、上記5種の検索を共通の処理で実現し、あわせて検索ファイルの小容量化による検索ファイルの主記憶常駐および検索の高速化も実現できる検索方式を述べる。

2. 検索方式

2.1 キーワードのコード化

KWを分野別など、情報を整理利用するのに適した登録番号を付加する。このKWを一文字ごとに分解し、各文字に先頭からの文字位置を示す文字番号を与える。この登録番号と文字番号を式1に代入し、文字コードを作成することでキーワードをコード化する。例えば、「ISDN」に登録番号125を与える場合、文字コードは図1のようになる。

$$\text{文字コード} = \text{登録番号} \times 100 + \text{文字番号} \dots \text{式1}$$

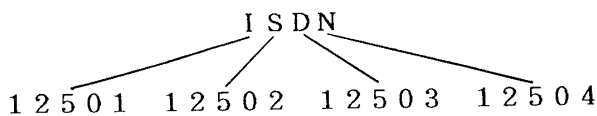


図1 KWコード化の例

文字コードを4バイトで構成すると、最大100文字長のKWを、 $2^{32}/100 = \text{約}4\text{千万種}$ 、整数として取り扱うことができる。また、KW長が最大20文字の場合には、 $2^{32}/20 = \text{約}2\text{億種}$ のKWを扱うことができる。

2.2 検索ファイル

検索ファイルは、システムで使用可能な各文字に対応するブロックで構成される。例えば、JIS第1及び第2水準漢字表(以下JIS表と略す)の各文字をUNIXコードで取り扱う場合は、式2により変換したコードをブロック名としている。KWを登録する場合には、KWの各文字に対応するブロック内の末尾に、該当する文字コードを格納する。「ISDN」を登録する例を図2に示す。このとき、各ブロック内の文字コードは、昇順に格納されることになる。

$$\text{ブロック名} = \text{UNIXコード} - A1A0 \dots \text{式2}$$

D	----	1 2 5 0 3	-----
I	-----	1 2 5 0 1	-----
N	---	1 2 5 0 4	-----
S	-----	1 2 5 0 2	-----

図2 検索ファイル登録例

文字コードを4バイトとすると、ファイル容量は  $4\text{バイト} \times \sum_{i=1}^n (\text{KW文字数})_i$  で小規模になるので主記憶常駐が可能になり、検索の高速化も実現できる。

2.3 検索方法

検索時は、入力された文字列を一文字毎に分解し、式2からブロック名を算出し、文字順に該当するブロックを検索ファイルから読み出す。次いで式3に示すように、ブロック間で登録番号が等しく文字番号が1差の文字コードを取り出す。

(文字順nのブロック内文字コード) -  
 (文字順n+1のブロック内文字コード)  
 = -1・・・式3

この後、検索指示に従い、前方一致/後方一致/完全一致/部分一致/中間一致特有の処理を行い、文字コードから検索対象の登録番号を取り出す。

例えば、検索入力ISで前方一致検索を行う場合、図2からIのブロックとSのブロックを取り出し、式3を使用して該当する文字コードを抽出すると、この中に文字コード12502も含まれる。次に、前方一致特有処理として、検索入力ISの文字数2に等しい文字番号を有する文字コードを捜すと、この中に12502も含まれ、「ISDN」が検索されていることがわかる。残り4種の検索も同様に文字番号の処理を変更することで実現できる。

#### 2. 4 その他の処理

追加登録は、2. 1と2. 2で説明した登録処理同様に、KWの各文字に該当するブロックの末尾に新規文字コードを追加することで行う。また、削除は、削除KWの各文字に該当するブロックの該当文字コードを特殊記号に変更し、登録番号を欠番にする。この結果、追加登録と削除を短時間でできるようになった。

### 3. 実験結果

標記方式を評価するために、PFU社製A-50ミニコン(OS:UNIX)上に、51000冊の図書名を登録した実験システムを構築した。洋書の検索も考慮し、KW長を最大100文字として式1で文字コードを作成し、実験では、検索ファイル容量、検索時間、追加登録時間、削除時間を調べた。ファイル容量は約1.8Mバイトで、各処理時間は表1にとおりである。なお、使用した図書名は1文字~22文字で、平均は9文字である。また、検索時間は検索ファイルから該当する全登録番号を取り出すまでの処理時間、追加登録時間はディスクと主記憶の検索ファイルにKWを登録するまでの処理時間、削除時間はディスクと主記憶の検索ファイル内の登録データを特殊記号に変更するまでの時間である。

	最小	最大	平均
完全一致検索	4 ms	1 2 1 ms	7 0 ms
前方一致検索	2 ms	1 7 5 ms	3 2 ms
後方一致検索	3 ms	3 4 5 ms	3 2 ms
部分一致検索	3 ms	3 2 2 ms	3 1 ms
中間一致検索	4 ms	7 8 0 ms	3 2 ms
追加登録	3 8 S	3 9 S	3 8 S
削除	3 5 S	4 8 S	3 6 s

表1 処理時間

#### 4 おわりに

KWをコード化して取り扱うことにより、完全一致/前方一致/後方一致/部分一致/中間一致の各種検索を共通の検索ファイルで実現できること、検索時間が完全一致を除くと平均32msであること、また、検索ファイルがパソコンやワークステーション主記憶に展開できるほど小規模に構成できることを確認した。さらに、追加登録時間が平均38s、削除時間が平均36sであり、データベース・メンテナンス面も実用的な速度と判断される。

また、本方式は、KWをその登録番号と文字番号で管理できるが、これはKWだけではなく、情報の総称に登録番号を付加し、情報の構成部分に出現番号を与えることにより、情報一般の登録検索に応用できると考えられる。

さらに本方式は複数のブロックから構成される検索ファイルに、コード化された文字コードを格納する方式なのでハード化が容易であることから、図書や辞書検索以外の各種情報検索にも応用可能と思われる。なお、意図的に平仮名表現がされている、あるいは、複数の読みが可能なKWを対象とする検索への適用が今後の課題である。

#### 5 謝辞

本テーマの取り組みにあたって助言をいただいた当研究所第一研究部外月前部長(現東京電力システム研究所主席研究員)、日頃から御指導いただいている第一研究部藍沢部長および有益な議論をいただいた山田主任研究員はじめ第一研究部各位に感謝いたします。

#### 参考文献

- [1] 福島・菊地、文字列検索LSIを用いた国語辞書システムの構築法、情処37全大5B-10、1988