

# 高速先頭照合方式による ストリングサーチ高速化の検討

7H-6

秋沢 充           野口孝樹           川口久光  
加藤寛次        畠山 敦           志村隆則  
(株)日立製作所           中央研究所

## 1. はじめに

情報処理システムの記憶装置の容量が増大するに従い、文書データを漏れなく高速に検索する処理が一層重要となっている。しかも、一般ユーザが簡単に文書データの蓄積・検索をしたいという要求がある。このため、自由に設定したキーワードによる高速の全文検索の必要性が高まっている[1][2][3]。

全文検索において重要なストリングサーチアルゴリズムのひとつとして、複数キーワードを一括して検索する有限状態オートマトンを用いた方法[4][5] (以下FSA法と記す) が知られている。本報告ではFSA法をベースとしたストリングサーチの高速化方式について述べる。

## 2. FSA法高速化の課題

図1にFSA法の処理の概念図を示す。一般に状態遷移テーブルはメモリで構成され、その容量は大きく、状態遷移を制御する部分とは別チップ構成となる。文書データが1文字入力されるたびに状態遷移テーブルをアクセスし、状態遷移のシーケンスが繰り返される。したがって状態遷移テーブルへのアクセスが毎サイクル必要となり、処理速度向上の妨げとなる。

今回我々は、状態遷移テーブルのアクセス頻度を減らすことで高速化が可能であることに着目し、新たな高速ストリングサーチ方式を検討した。

## 3. 高速先頭照合方式の提案

本報告で提案する高速先頭照合方式は、各状態への遷移頻度の差に着目した方式である。ハードウェア論理によって遷移頻度の高い状態におけるテーブルアクセスを不要とする。これによりストリングサーチを高速化する。

図2にFSA法の解析モデルを示す。簡単のためにキーワードが“ $C_0.C_1 \dots C_k \dots C_{n-1}$ ” ( $C_k$ :任意の文字) のみの場合について考察する。レベルkにおいて状態遷移処理が行なわれる場合には、必ずレベル0、

…レベル(k-1)での処理が行なわれている。したがって、従来のFSA法ではレベル0が状態遷移テーブルのアクセス頻度が最も高く、順次レベル1, レベル2, …と低くなる。高速先頭照合方式ではレベルk = kcまでをハードウェアによって高速に処理(先頭照合処理)し、レベル(kc+1)以降は有限状態オートマトンの状態遷移により処理(後方照合処理)するものとする。

[定義1] 任意の状態での状態遷移処理において、これをソフトウェアで実行した場合のステップ数に換算したものを処理量と定義する。

[定義2] レベル0からレベルkまでの処理量の総和の、レベル0からレベル(n-1)までの全処理量に対する割合WR(k)を相対部分処理量と定義する。

ある特定の文字(字種数Nc)が文書データ中出现する確率を1/Ncであるとする、これは状態遷移確率fに等しいと考えることができる。このとき、相対部分処理量WR(k)は次式となる。

$$WR(K) = \frac{1 - f^{k+1}}{1 - f^n} \quad (1)$$

( $0 \leq k \leq n-1$ , ただし  $f = 1/Nc$ )

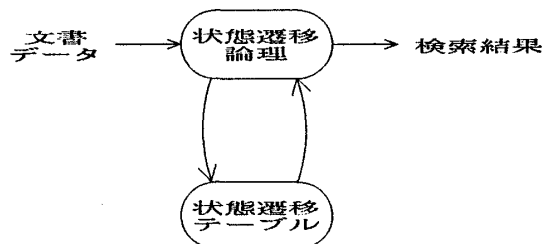


図1 FSA法の処理概念図

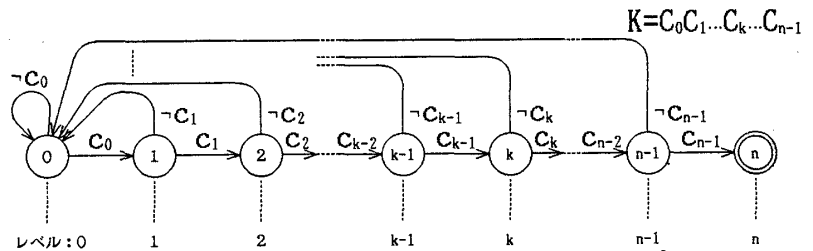


図2 有限状態オートマトンの解析モデル

Proposal of New String Search Method with Fast Partial Pattern Matching

Mitsuru Akizawa, Kouki Noguchi, Hisamitsu Kawaguchi,  
Kanji Kato, Atsusi Hatakeyama, Takanori Shimura  
HITACHI, Ltd.

“a b c d e”を検索する場合のWR(k)を、一例として求めた結果を図3に示す。kc=2とすると全処理量の99%を、kc=1としても全処理量の96%を先頭照合で処理できる。以上の評価から、キーワードの先頭の2~3文字の比較処理が全処理量の大部分を占めることがわかる。したがって、この部分をハードウェア化して高速に先頭照合処理をすることで、ストリングサーチが高速化される。

図4はkc=1の場合に、本方式を並列比較器と後方照合用に生成した有限状態オートマトンとで構成した概念図である。長方形で囲まれた2文字を並列比較器へ設定する。並列比較器からの一致信号がオートマトンの初期状態を発火させ、以後は状態遷移テーブルに従って文書データとの比較照合処理を行っていく。

4. 予測性能の評価

高速先頭照合処理方式を採用したストリングサーチの加速率の評価を行なう。ストリングサーチ処理全体の加速率Atは次式より求められる。

$$1 / A_t = (R_c / A_f) + (1 - R_c) \quad (2)$$

$$\text{ただし、} R_c \equiv W R(k_c) \quad (3)$$

ソフトウェア処理を基準とした、並列比較器による先頭照合処理の加速率Afを80、日本語文書における相対先頭照合処理量Rcを96%(kc=1:先頭2文字)と仮定する。これらの仮定から相対ストリングサーチ処理時間(1/At)を算出した。

算出結果から、日本語文書のストリングサーチにおいて、高速先頭照合方式(kc=1:先頭2文字)により約20倍の加速率が達成可能であることがわかった。

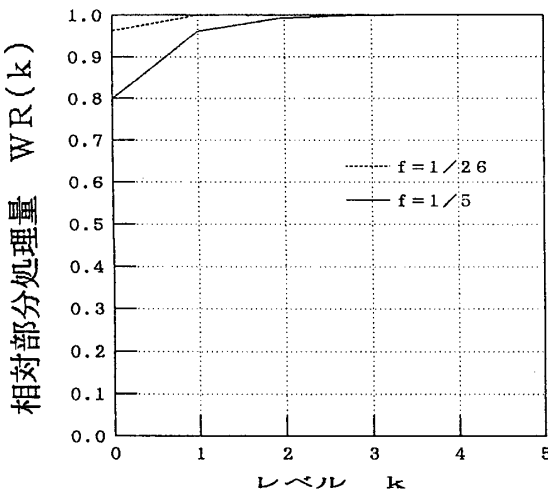


図3 解析モデルによる相対部分処理量

5. おわりに

FSA法をベースとしたストリングサーチを高速化する高速先頭照合方式を提案した。従来のソフトウェア処理に比較して、約20倍の加速率となる見通しを得た。

参考文献

[1]加藤,他:大規模文書情報システム用テキストサーチマシンの研究,情処研報,Vol.89, No.66, pp.14.6.1-14.6.8(1989.7)  
 [2]高橋,他:フルテキストサーチのハードウェア技術について,同上, pp.14.5.1-14.5.8(1989.7)  
 [3]早川,他:ストリームデータプロセッサSDP(1)(2),情処第37回全大, pp.113-116,(1988)  
 [4]A.V.Aho, M.J.Corasick:Efficient String Matching, Commun. ACM, Vol.18, No.6, pp.333-340(1975.6),  
 [5]篠原,有川:日本語テキスト用のAho-Corasick型パターン照合アルゴリズム,情処研報, Vol.86, No.48, pp.52.4.1-52.4.8(1985.11)

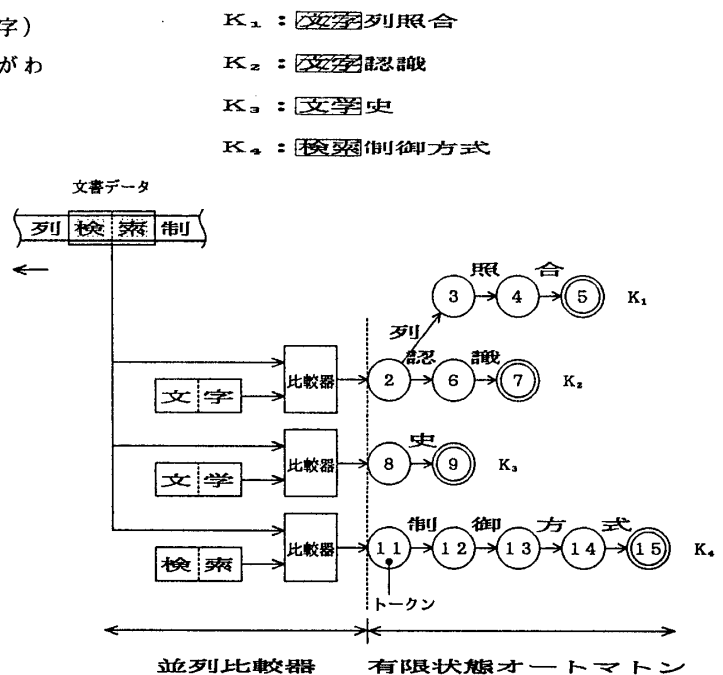


図4 並列比較器とFSAで構成した高速先頭照合方式の概念図