

ストリームデータプロセッサ SDP (2)

パターンマッチング方式

3N-10

伊藤 正雄 早川 佳宏 田村 登 上田 謙一
松下電器産業株式会社 東京研究所

1. はじめに

本稿では、大量のデータベースから検索、更新等の処理を高速に行うストリームデータプロセッサSDPのパターンマッチング方式とハードウェア構成について報告する。

パターンマッチングのハードウェアに適したアルゴリズムとしては有限状態オートマトン方式[1, 2]、連想メモリ方式[3]、セルラレイ方式等があるが、本方式はストリームデータに適し、かつ柔軟なパターンマッチングが可能な有限状態オートマトン方式を用いる。有限状態オートマトン方式は、文字列を有限の入力文字と状態の2次元からなる遷移関数で表現し、入力文字で遷移関数を参照し状態を変化させることにより検索処理を実現するものである。有限状態オートマトンの利点は、正規表現を用いた複雑なパターンの検索が可能で、遷移関数をメモリで実現するとハードウェアが簡単で済むことである。しかし入力文字の種類が多いと遷移関数を格納する領域が非常に大きくなり、かつ前処理時間がかかるという欠点がある。

SDPは今後1チップ化を考えており、有限状態オートマトンをそのままインプリメントしたのでは、膨大な作業領域がボトルネックになると考えられる。よって本手法では、作業領域を縮小させ、かつ有限状態オートマトンの特徴を損なわない方式を考案し、インプリメントを行った。

2. SDPの検索機能

SDPの検索機能としては自然語検索ができるような柔軟な検索機能を第一目標とし、UNIXのlexと同等の正規表現を用いた検索を実現した。またデータベースをレコードの集合、レコードをフィールドの集合とした場合、レコード全体に対して、また指定したフィールドに対して検索を可能にした。

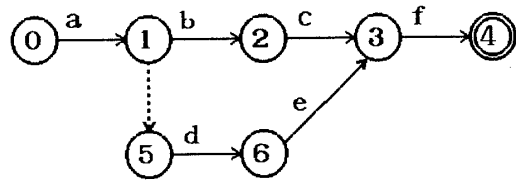
SDPの主な検索機能を以下に示す。

- ・フィールド指定検索
- ・正規表現を含む検索
- ・オーダ検索
- ・論理演算子(AND, OR)を用いた複合条件検索

3. パターンマッチング方式

SDPのパターンマッチング方式は先に述べたように、有限状態オートマトン方式を用いている。有限状態オートマトン方式で有名なものとしてAC法[4]がある。AC法は、goto、failure、outputの3種類の関数を用いて検索を行うものであり、データを一回走査するだけで、複数個のパターンを検索できることが特徴である。しかし

AC法でもgoto関数という2次元の遷移関数を用いているため膨大な作業領域が必要である。本手法では、そのような2次元の関数を用いない。その代わりに、比較が失敗した場合に遷移する状態を示すfailure関数と、比較が成功した場合に遷移する状態を示すsuccess関数と、検索パターンを格納するpattern関数と、検索の成功、失敗を示すoutput関数を用いる。図1に遷移図の例を示す。実線がsuccess関数、点線がfailure関数、実線の上の文字がpattern関数、二重円が最終状態を示す。



a (b c | d e) f
図1 遷移図の例

以上の関数だけでは、可変長文字列が含まれる場合、文字範囲を指定するような検索の場合、効率が悪くなる欠点がある。そこで表1に示すように以下の関数を加える。

・input_control関数

比較が成功、失敗した場合に入力データのポインタを1つ進める、または進めないことにより、柔軟な比較操作を可能にするものである。

・compare関数

数値、文字の大小比較をする場合に、比較の種類(=、>、<、≥、≤、≠)を指定する。

また、特定のフィールドに対して検索を行うために、フィールドの設定を示すfield_set_flag、フィールド番号を指定するfield_number、フィールドのデリミタを検出した場合に遷移する状態を示すlink関数、比較するbyte幅を柔軟に変更できるように、byte_widthを用いる。

以上挙げた関数(まとめて検索テーブルと呼ぶ)は全て1次元の配列であり、かつbit数は大きくないためコンパクトな作業領域で済む。表1にa (b | c) dを条件とした検索テーブルの例を示す。表1でBW, FNはFSが1の場合に有効になり、OP, CO, ICはFSが0の場合に有効になる。

4. ハードウェア構成

図2にSDPパターンマッチ部(PME)のハードウェア構成を示す。検索するデータベースは入力時にレコード単位でバッファリングされる。RECORD_BUFFERは入力、

Stream Data Processor SDP (2)

The Method of Pattern Matching

Masao Ito, Yoshihiro Hayakawa, Noboru Tamura, Kenichi Ueda

Tokyo Research Laboratory Matsushita Electric Industrial Co., Ltd.

表1 検索テーブルの例

状態	PA	FA	SU	LI	FN		FS
					OP	IC	
0	-	-	-	-	01	1	1
1	a	6	2	6	00	0	0
2	b	3	4	6	00	0	0
3	c	6	4	6	00	0	0
4	d	6	5	6	00	0	0
5	-	-	-	-	01	-	0
6	-	-	-	-	10	-	0

PA :pattern OP :output
 FA :fail FN :field_number
 SU :success CO :compare
 LI :link IC :input_control
 BW :byte_width FS :field_set_flag

パターンマッチ、出力のそれぞれがパイプラインで処理できるように3バッファ方式を採っている。また入力時に、フィールド単位の先頭アドレスをRECORD_BUFFERに格納することにより、検索時に指定したフィールドを直接読み出すことを可能にした。

PMEのハードウェア構成を実行手順を示しながら説明する。前処理として検索テーブルをPATTERN_BUFFERに転送する。次にPMEの起動がコントローラからかかるとPATTERN_BUFFERの先頭から読み出す。field_set_flagが有効の場合は、field_numberでフィールドの先頭アドレスを得、それを用いて更にRECORD_BUFFERからデータを読み取る。次の状態でPATTERN_BUFFERのパターンとRECORD_BUFFERから読み取ったデータと比較を行う。本処理は、32bit単位で比較を行うことが可能な完全一致検索と、文字コードのbit数でしか比較ができない可変長文字列等の正規表現を用いた検索を同一のハードウェアで行っている。これを実現するために8bit単位で比較を行う場合

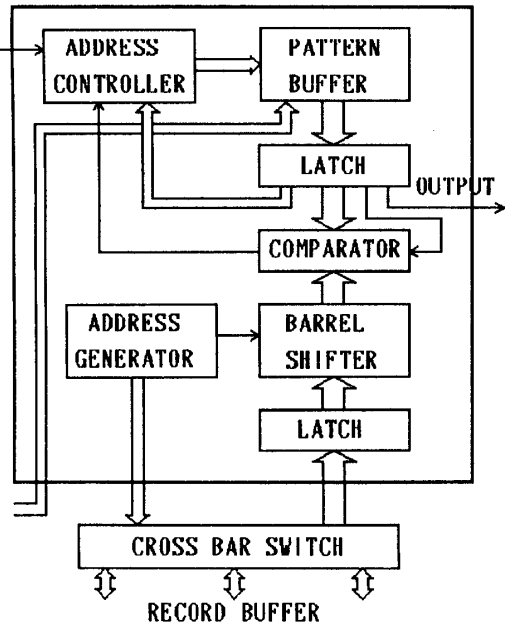


図2 PMEのハードウェア構成

には、BARREL_SHIFTERを用いてPATTERN_BUFFERのパターンと同じbit位置にシフトを行なう。比較器の結果はADDRESS_CONTROLLERに入力され、結果によりsuccess、fail、linkが示す状態に遷移する。

5. 性能予測

ここではPMEの性能予測について述べる。SDPは、フィールド単位の先頭アドレスをインデックスとして持っているため、複数フィールドで1フィールドに対する検索は速いと考えられる。ここでは全文テキストデータのように全データから検索する場合について評価する。またSDPは主記憶装置として大容量RAMを対象としているため、2次記憶系についてはここでは考えないものとする。比較する対象としては、汎用CPUを用いて有限状態オートマトン方式で検索する場合とする。以下に性能予測環境を示す。

- CPU 68020、16MHz、約1MIPS
- メモリアクセス時間（アドレス変換、バスの遅れを含む）450ns
- 使用データサイズ 2Mbyte

以下に汎用CPUの処理時間Tcを示す。

$$T_c = D_s \times B_s \div M_i$$

Dsはデータベースサイズ

Bsは1byte処理あたりのステップ数

MiはCPUのMIPS値

Bsを約15とするとTcは30秒となる。

次にSDPの処理時間Tsを示す。

$$T_s = D_s \times T_b \times \alpha$$

Tbは1byteあたりの処理時間

αは属性読み出しとフィールド設定時間

TbはPMEのサイクル時間で200ns、αを約1.25とするとTsは0.5秒となり汎用CPUより2桁処理速度が速くなることが予想される。

また作業領域の観点から見ると、8bitの文字で最大状態数256では本方式は2Kbyteで済むが、有限状態オートマトンでは64Kbyte必要となるため、32分の1の領域で済む。

6. おわりに

ストリームデータプロセッサのパターンマッチング部として、AC法のgoto関数を用いず、failuer関数を用いて作業領域を小さくし、かつ複雑な検索を可能にした検索方式とハードウェア構成について述べた。今後、実際のアプリケーションで評価を行う予定である。

【参考文献】

[1]R. W. Floyd, J. D. Ullman: "The Compilation of Regular Expressions into Integrated Circuits", Journal of ACM, Vol. 29, No. 3, July (1982)
 [2]井上他: "情報提供サービスに適用可能な超大規模リレーショナル・データベースマシン", データベースシステム, 47-5 (1985)
 [3]Takahashi, k.: "A New String Search Hardware Architecture for VLSI", Proc 13th ISCA, pp.20-27, June, (1986)
 [4]Aho, A. V., M. J. Corasick: "Efficient String Matching: An Aid to Bibliographic Search", CACM Vol. 18, No. 6, (1975)