

BM法を用いたハードウェア方式による 複数文字列を同時検索するアルゴリズムの提案

長田 千加子、大曾根 匡
専修大学 経営学部 情報管理学科

1. はじめに

【1】では、BM法をハードウェアを用いて並列化した文字列検索アルゴリズムを提案した。本論文では、それを複数パターンの同時検索を可能にするように拡張する。

2. ハードウェア構成

単数パターンに対する検索から複数パターンに対する検索に拡張するにあたって、ハードウェア構成の変更した点を以下に示す。ここで、説明を容易にするため、パターン数 n は2とする。

- ①パターンレジスタの個数を2個にする (P_1 と P_2)。
- ②データレジスタ D とパターンレジスタ P_1 と P_2 の文字同士を同時比較できるように、並列比較ユニット PCU も2個にする。
- ③比較結果ラッチも C_1 と C_2 の2個用意し、比較結果を別々に記憶する。

3. 動作例

例として、パターン1が「ABCD」、パターン2が「EFGH」の場合について考える。また、テキストは図2に示す。

この例の場合、あらかじめパターンレジスタ P_1 に「ABCD」、 P_2 に「EFGH」を格納しておく。初期縦読みアドレスは3に設定しておく。

(1) 時刻1：縦読みアドレス3から最初の縦読みデータ「IMPS」をフェッチし、データレジスタ D に取り込む。そして並列比較ユニット PCU_1 と PCU_2 を用いてデータレジスタ D とパターンレジスタ P_1 、 P_2 の各文字の比較を行う。この場合、両者の文字同士で一致する文字はないので、比較ラッチ C_1 と C_2 は全て0となり、横読みを行わないことを判断する。そして、縦読みアドレスを+16更新し19とする。

(2) 時刻2：縦読みアドレス19から2回目の縦読みを行い、データ「JCNF」をデータレジスタ D に取り込む。このとき、データレジスタ D

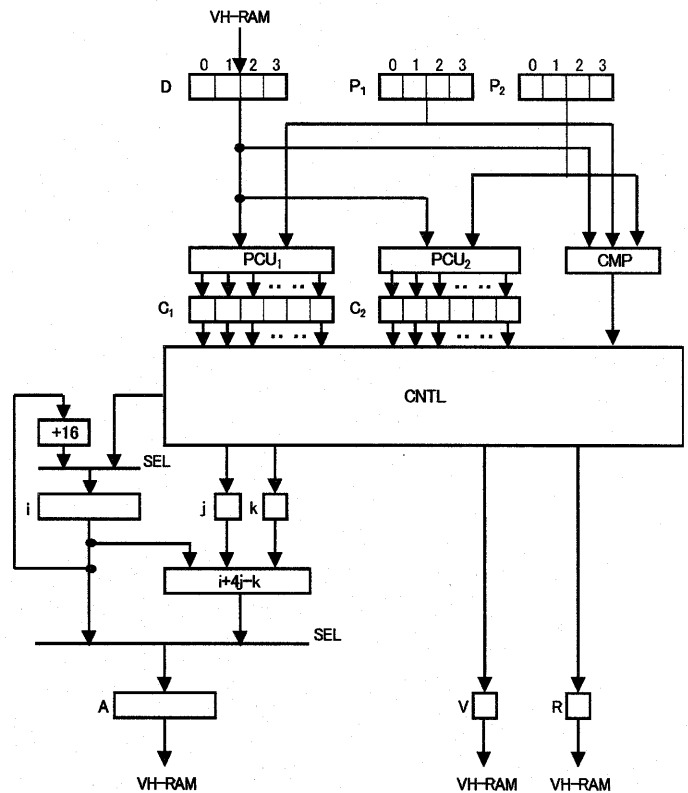


図1 検索エンジンのハードウェア構成

の1文字目「C」とパターンレジスタ P_1 の2文字目の「C」が一致するので、ラッチ $C_1[1, 2]$ が1となる。そこで、横読みアドレス21を算出する。同様に、ラッチ $C_2[3, 1]$ も1となる。これより、横読みアドレスは30となる。

		テキスト			
0	A	B	C	I	
4	F	A	B	M	
8	C	D	A	P	
12	B	C	D	S	
16	C	A	B	J	
20	D	A	B	C	
24	D	D	A	N	
28	A	B	E	F	
32	G	H	A	V	
36	B	C	D	L	
40	U	R	B	B	
44	D	B	K	Q	

(3) 時刻3：横読みアドレス21からテキスト「ABCD」

図2 テキストの例

を横読みし、データレジスタ D に格納する。この場合、パターンレジスタ P_1 とデータレジスタ D の内容が一致するのでパターン1を検出する。

(4) 時刻4：同様に、横読みアドレス30からテキスト「EFGH」を横読みし、データレジスタ D に取り込む。この場合もパターンレジスタ P_2 とデータレジスタ D の内容が一致しているので、パターン2を検出する。そして、すべての横読みが終わったので、縦読みアドレスを35に更新し、縦読みを再開する。

(5) 時刻5：縦読みアドレス35からデータ「VLBQ」をデータレジスタ D に取り込む。このとき、ラッチ $C_1[2, 1]$ が1となる。そこで、パターンレジスタ P_1 に対して、横読みアドレス42を算出する。

(6) 時刻6：横読みアドレス42からテキスト「BBDB」を横読みし、データレジスタ D に格納する。そして、パターンレジスタ P_1 と比較する。この場合は一致しないので、パターン1がここに存在しないことがわかる。

このようにして、最終的には、縦読み3回、横読み3回の計6回のデータ読み出しで、2つのパターンを同時に検索することができる。

4. 性能評価

パターン数 n を複数にしたことによって、読み出し回数などがどのように変化するかをシミュレーション実験した。その結果を図3に示す。縦軸に読み出し回数を、横軸にアルファベットの文字の種類数を示している。

アルファベットの文字の種類数は、4、8、16、32、64、128と変化させた。テキスト長は、10000文字とし、ランダムに発生させた。また、パターンについては、長さを4文字とし、ランダムに生成した。パターン数 n は、1から4まで変化させた。

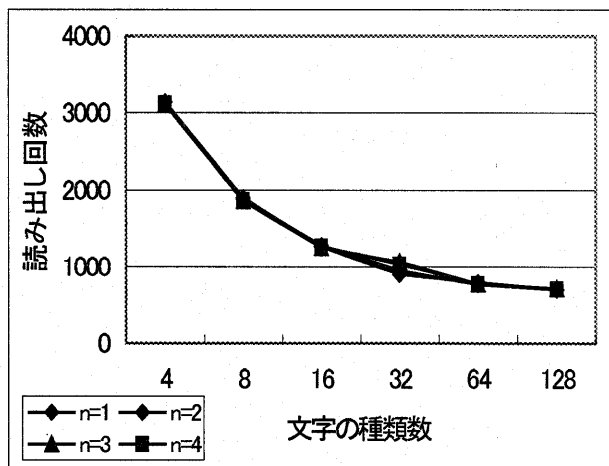


図3 読み出し回数の変化

この図から、文字種類数が増加することによって、読み出し回数が減少していることがわかる。そのことより、文字の種類数が増加するほど、アルゴリズムの性能が良くなることがわかる。

また、パターン数 n が増加しても、読み出し回数は大きくは変化しないことがわかった。

参考文献

- [1] 大曾根匡：「BM法を用いたハードウェア方式による文字列検索アルゴリズムの提案」、IPSJ第61回全国大会 1Q-02、2000.10.