

複数文字列の探索に関する研究

浦谷 則好

①

# 複数文字列の探索に関する研究

平成8年9月

浦谷 則好

# 目次

第1章 序論	1
1.1 本研究の背景と目的	1
1.2 文字列検索の手法	2
1.2.1 インデックス検索手法	2
1.2.2 長尾・森の方法	4
1.2.3 全文検索手法	6
1.3 本論文の構成	8
第2章 複数文字列の探索手法	9
2.1 まえがき	9
2.2 全文検索手法	9
2.2.1 素朴な方法	10
2.2.2 Knuth-Morris-Prattの方法	10
2.2.3 Boyer-Mooreの方法	11
2.2.4 Aho-Corasickの方法	14
2.3 FAST法のアルゴリズム	15
2.3.1 FAST法の考え方	15
2.3.2 FAST法のアルゴリズム	18
2.4 あとがき	25
第3章 FAST法の性質	27
3.1 まえがき	27
3.2 FAST法の効率	27
3.3 正当性の証明	32
3.3.1 完全性の証明	33
3.4 検索実験	38
3.5 メモリ削減手段	46
3.6 あとがき	49

第4章 FAST法の2次元パターン検索への拡張	51
4.1 まえがき	51
4.2 FAST法の2次元パターン検索への拡張	51
4.2.1 2次元FAST法のアルゴリズム	56
4.2.2 検索実験	61
4.3 FAST法のn次元パターン検索への拡張	63
4.4 あとがき	64
第5章 結論	65
謝辞	67
参考文献	68
本論文に関する対外発表	68
本論文関連以外の主な対外発表	69

## 第1章 序論

### 1. 1 本研究の背景と目的

情報化社会の進展に伴って、個人個人に供給される情報はどんどん増大し、24時間を費やしても、もはやその全てにアクセスすることができなくなっていると言われている。必要な情報を、必要な時に利用できるようにデータベースが様々な機関で作成され、供給されている。画像や音声情報のデータベース化も進んできているが、未だに文字情報が主流である。データの検索の手掛かりとしては、キーワードが最も利用されている。

キーワード検索手法としては、キーワードに対する索引表をあらかじめ構築しておくインデックス検索手法<sup>1)</sup>が一般的であるが、この手法には次のような問題点がある。

- (1) キーワードの抽出に多大な労力と時間を必要とする。
- (2) 検索対象はキーワードにより制約され、任意の語彙の検索は不可能である。

そこで、あらかじめデータに何もキーワードを付けておかないで、文書中にある語を全てキーワードとして使えるようにする全文検索が考えられている。この場合、ユーザは何の制約も受けずに自由にキーワードを指定できることになる。データベースシステムはあらかじめユーザの要求を知ることができないので、検索要求が出る度にユーザが指定したキーワードの出現有無をデータに直接当たって調べなければならない。したがって、データベースが大きくなるに連れて、検索時間もどんどん大きくなっていく。そこで、できるだけ高速な文字列検索手法が求められている。

また、ワードプロセッサやプログラミング言語のコンパイラ（やインタープリタ）においても文字列の検索は基本的な機能として必須のものである。文字列の検索は、指定した文字列（以下パターンと呼ぶ）が文書（以下テキストと呼ぶ）中のどこに現われているかを調べるという単純な作業であるが、前述のように電子化された文書の処理にとっては全体の効率に大きく関与する操作なので、これを効率的に実現するためのアルゴリズムは各所で研究されている<sup>2)</sup>。

1個のパターンを探す手法としてはBoyer-Mooreの方法（BM法）<sup>3)</sup>が最も効率の良い方法として知られている。しかし、この方法で複数パターンの検索を行うためにはパターンの数だけの繰り返し使用が必要となる。複数のパターンを同時に効率的

く検索する方法としては、Aho-Corasickの方法（AC法）<sup>4)</sup>が有名である。このprobe rate（テキスト1文字当たりの照合回数；つまり効率の逆数）は常に1になる。

本論文ではBM法にAC法の状態遷移という考え方を取り入れて、BM法を複数パターンに拡張する方法を研究する。BM法を複数パターンに拡張しようとする試みとしてはKowalskiとMeltzerによる方法<sup>5)</sup>やCommentz-Walterによる方法<sup>6)</sup>があるがここで紹介する方法はそれよりprobe rateが優っている。

この方法の効率を推定し、実験によりそれが正しいことを示す。また、アルゴリズムの正当性を証明する。さらに、この方法を2次元パターンの検索に拡張する方法についても明らかにする。

## 1. 2 文字列検索の手法

### 1. 2. 1 インデックス検索手法

キーワード検索手法としては、キーワードに対する索引表をあらかじめ構築しておくインデックス検索手法が一般的である。この手法の概念図を図1. 1に示す。つまり、あらかじめキーワードを人手あるいは自動的に抽出しておいて、それが出現する位置をインデックスファイルに記憶しておく方法である。キーワードの検索時には、索引表からキーワードを検索し、キーワードの位置情報を得ればよい。したがって、非常に高速な検索が実現できるが、

- (1) 索引表をとっておくためにかなり大きな記憶容量を必要とする。
- (2) 新たなキーワードを追加するには索引表を作り直す必要があり、この時多大な時間を必要とする。
- (3) キーワードの自動抽出法も研究されている<sup>7),8)</sup>が、キーワードの抽出に多大な労力と時間を必要とする。

という問題点があり、何よりも

- (4) 検索対象はキーワードにより制約され、任意の語彙の検索は不可能である。といった致命的な欠点を有している。



図 1. 1 インデックス検索の概念図

### 1. 2. 2 長尾・森の方法

インデックス検索手法には「検索対象はキーワードにより制約され、任意の語彙の検索は不可能」という致命的な欠点があった。それなら、「テキストに現われる全ての文字列に対してインデックスを作成しておけばよい」という考えもありうる。しかし、例えば  $n$  文字列の全てのインデックスを作成するには、文字列を記憶するのに約  $2nk$  バイト ( $k$  はテキストの長さ；日本語は2バイトで表されるから2が掛かる)、文字の出現位置を記憶するのに  $4k$  バイト (ポインタを記憶するのに4バイト必要と仮定して)、合わせて  $2k(n+2)$  バイトの記憶容量が必要である。したがって、長さ  $n$  までの文字列全てのインデックスを作成するには  $(n^2+5n)k$  の記憶容量が必要だということになって、 $n$  が3でも元のテキストの12倍、 $n$  が10なら元のテキストの75倍もの記憶容量が必要となる。したがって、元のテキストが大きければインデックス検索手法によって任意の語彙を検索するということは事実上不可能となる。

長尾と森はテキストから  $n$  グラム統計を取る際に、記憶容量をそれほど必要としない方法を考案している<sup>9)</sup>。この手法は文字列検索のために考案されたものではないが文字列検索にも利用することができる。

この手法では、テキストの  $i$  文字目から最後の文字 ( $k$  文字目) まで (実際には最大255文字分としている) を1つの文字列と見なし、 $i=1$  から  $i=k$  までの  $k$  個の文字列の存在を仮定する (実際には図1、2の原ストリングのように1本の文字列である)。この  $k$  個の文字列を辞書順にソートする。ソートされた結果は文字列の配列順序によって示される。その簡単な例を図1、2に示す。各文字の位置を示すために一般に4バイトを必要とするので、このソートされたデータの記憶容量は  $4k$  バイトである。

ポインタ	1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
原ストリング	ある日ある場所で会合がもたれた。
ソート結果	16 1 4 11 8 15 13 12 2 5 14 3 9 10 7 6

図1. 2 文字列のソート例



このようにソートされた文字列に対して、隣接している2つの文字列が先頭から何文字目まで一致しているかを調べ、これを前の文字列の情報として図1. 3のように付加する。一致文字数を最大255までとすればこの情報を保持しておくのに1文字列当たり1バイトが必要なので全体では $k$ バイトが必要である。以上の処理に必要な記憶容量は、(1)長さ $k$ の原テキストの記憶に $2k$ バイト、(2)ソートされた結果の記憶に $4k$ バイト、(3)一致文字数の保持のために $k$ バイト、となり全体で $7k$ バイト必要であることがわかる。言い換えれば、付加すべき記憶量は元のテキストの2.5倍で良いということになる。これは、インデックス検索手法よりずっと少ない量である。

今、ある文字列をテキストから検索する場合を考える。この文字列がテキストに存在するか否かはソート結果があるからバイナリーサーチでポインターの示す文字列と比較することで簡単に調べることができる。これに要する時間は $o(\log_2 k)$ で済むことは良く知られている。また、出現位置が1つでない場合も、バイナリーサーチで最初に見つけた位置から前後に一致文字数を調べていだけで文字列そのものを比較することなく、どれだけを検索出力があるのかを知ることができる。この検索時間はもちろん直接検索結果の得られるインデックス検索手法と比べれば大きい、 $o(\log_2 k)$ であることを考えればかなり小さいと言える。

部分文字列	一致文字数
ある日ある場所で・・・	5
ある日ある人が・・・	2
ある場所で会合が・・・	6
ある場所で会う。	2
あるのは・・・・・・・・	1
あればあったで・・・	0
いい加減な・・・・・・・・	

図1. 3 一致文字数の情報

### 1. 2. 3 全文検索手法

インデックス検索手法では任意の語をキーワードとして検索することができないことが分かった。長尾・森の方法では元のテキストに対して2.5倍もの余分な記憶容量が必要な上に、文字列のソートのためにかなりの前処理時間が必要である。

これらの欠点を克服するためには、あらかじめデータに何もキーワードを付けておかないで、文書中にある語を全てキーワードとして使えるようにする全文検索以外に手段がない。ユーザは何の制約も受けずに自由にキーワードを指定できることになるが、データベースシステムは予めユーザの要求を知ることができないので、検索要求が出る度にユーザが指定したキーワードの出現有無をデータに直接当たって調べなければならない。

全文検索法は、全文書を対象にそのまま語を検索する方法（全走査法）と文書をブロックに分けそのブロックの特徴ベクトル（重ね合わせ符号）を事前に照合することで、一部の文書ブロックのみを検索対象とする方法（部分走査法）<sup>(10),(11)</sup>に分けることができる。

部分走査法概念図を図1.2に示す。この手法では文書をブロックに分け、それに含まれる語彙を特徴ベクトル（文書ベクトル）で表しておく。検索したいキー

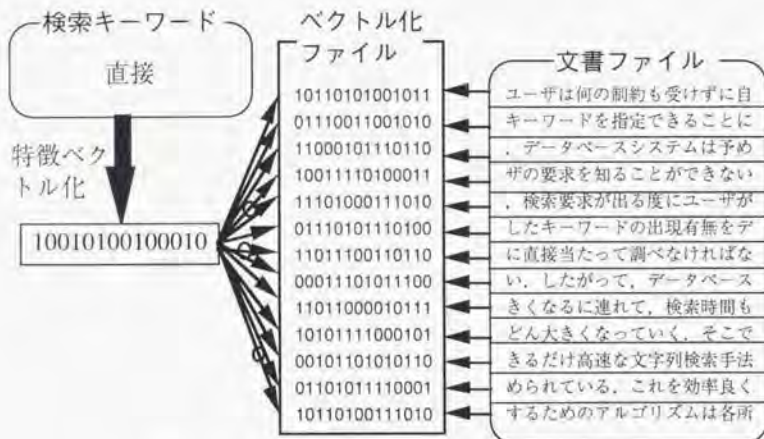


図1.4 部分走査法による検索の概念図

ワードも特徴ベクトルで表し、文書ベクトルとの論理積をとり、それがキーワードの特徴ベクトルと一致すれば、その文書ブロックを対象にキーワードの検索を行なう方法である。この方法の利点はうまく特徴ベクトルを作れば、検索対象とするブロックを絞ることができて高速な検索ができることである。しかし、ベクトル照合に成功しても、文書ブロックにはキーワードが存在しない場合もあり（これをFalse Dropという）、無駄な検索を行なう場合もある。False Dropを少なくするには特徴ベクトルを大きくすればよいが、そうすると文書ベクトルをとっておく記憶容量が大きくなるだけでなく、ベクトル照合に要する時間も長くなるという問題が生じる。そんな訳で、期待するほど全走査法に比べて検索時間を短縮できていない。そこで、以降は全走査法に絞って議論を進める。また、全走査法で使われる手法は部分走査法で文書ブロックを絞った後、キーワードの検索を行なう場合に使われるものと、本質的に何ら変わりはないので、より一般的でもある。

全走査法ではハードウェアによる高速化も試みられている<sup>12),13)</sup>が、柔軟性を考えればやはりソフトウェアによる高速化は現在でも重要な研究課題である。全走査法のアルゴリズムについては第2章で詳述する。

### 1. 3 本論文の構成

まず第1章で、研究の背景と目的、特にキーワード検索手法について概観し、合わせて、本論文の構成についても述べる。

第2章では全文検索手法のアルゴリズムについて検討する。従来のアルゴリズムについて説明した後、AC法の状態遷移という考え方を取り入れて、BM法を複数パターンの検索に拡張する方法(FAST法)を提案する。

第3章ではFAST法の効率について検討し、アルゴリズムの正当性を示す。さらに、実験結果を示し、効率の推定式がよく一致すること、多パターン時にはパターン数の増加に伴って効率が単調に減少しないという現象について説明する。メモリ削減手段についてもふれる。

第4章では、2次元パターン探索の従来のアルゴリズムについて説明した後、FAST法を2次元パターン探索に拡張する方法について述べる。実験により、従来の手法より効率が良いことを示し、さらに一般の $n$ 次元パターンの検索に拡張する方法についても述べる。

最後に第5章で結論を述べる。

## 第2章 複数文字列の探索手法

### 2.1 まえがき

近年、文書の作成や編集は、ワードプロセッサを用いて行なうことが多い。その際、しばしば使用する機能に、文字列の検索や置換がある。文字列の検索は、与えられた文字列が文書のどこに現われているかを調べ、見つかったらその位置にカーソルを移動し、見つからない場合にはそれをユーザに通知するものである。文字列の置換は、文書中のある文字列を別の文字列で置き換えるものであるが、当然、このためには文字列の検索を必要とする。コンパイラやインタープリタにおいても予約語の処理に文字列の探索が必要である。このように、文書进行处理するためには文字列の検索は最も基本的な操作といえることができる。このため、これを効率良く実現するためのアルゴリズムは昔から各所で研究されている。

文書データベースをあらかじめ決められたキーワードで検索する場合には第1章で説明したインデックス検索手法を使うのが一般的であるが、全く自由なキーワードをユーザに許すためには、この章で説明する全文検索手法を使わざるを得ない。この場合、第1章で述べた部分走査法による高速化も考慮することができるが、全走査法の方が根源的であることは前述した通りである。そのうえ、文書データベース検索以外のワープロやコンパイラなどの文書処理においても、全走査法の全文検索手法がオーバーヘッドの計算時間や記憶容量の面から考えて、最適であるのは明らかである。

### 2.2 全文検索手法

提案する文字列検索の新しいアルゴリズムの理解を容易にするために、まず従来の全文検索手法（このうちの全走査法）を概説する。最初に素朴な方法について説明した後、KnuthとMorrisとPrattが考案した照合回数がテキストの長さ に比例するアルゴリズム（KMP法）について説明する。BoyerとMooreは「パターンの右端から逆向きにテキストとの照合する」ことで、KMP法より効率の良い方法を考案した。AhoとCorasickは有限オートマトンの考え方を取り入れることによってKMP法を複

数のパターンの同時照合に拡張する方法を考案した。2. 2. 3と2. 2. 4でこれらの方法について説明する。

### 2. 2. 1 素朴な方法

文字列の照合は、あるパターンがテキスト中に含まれるか否かを判定し、含まれているならば、その位置を見つけることである。パターンが1つの場合には、この問題を解く次のような素朴なアルゴリズムを簡単に思いつくことができる。

まず、パターンの左端の文字をテキストの左端の文字に合わせるようにパターンをテキストの上に重ねておいて、パターンとその真下にあるテキストの部分文字列が一致しているかどうかの照合を行う。つまり、パターンとテキストの対応する文字が一致しているかどうかを左端から右方向へ1文字ずつ調べていく方法である。もし、途中で照合に失敗（文字が不一致）すると、パターンをテキスト上で1文字分右へずらし、再び照合を開始する。もちろん、パターンの全ての文字について照合に成功すれば、そこでテキスト中にパターンが見つかったことになる。

全てのパターン出現位置を探す場合、このアルゴリズムによる全照合回数は大体テキストの文字数に比例すると考えられるが、最悪の場合（例えば、全てaからなるテキスト"aaa...a"を対象にパターン"aaaab"を探すような場合）には、パターンの長さ（文字数）をm、テキストの長さをnとすれば、m回の照合で照合位置を1文字分ずらせるだけであるから、nがmより十分大きければ全照合回数は約m n回になってしまう。

### 2. 2. 2 Knuth-Morris-Prattの方法<sup>1)</sup>

最悪の場合でも照合回数がテキストの長さに比例するようなアルゴリズム（KMP法）をKnuthとMorrisとPrattが考案している。これはどこまで照合に成功したかの情報を保持しておいて失敗したときにパターンの検出漏れの心配のない最大の長さ分パターンを右にずらす方法である。

図2. 1にこの方法による照合の様子を示す。この図で3度目の照合（↑の所）に失敗した後、3文字分パターンをずらしているのは、以下の理由による。

パターンのi番目の文字を $\pi[i]$ とし、テキストのj番目の文字を $\tau[j]$ とすれば、

$$\pi[1]=\pi[3]=\pi[4]=a \neq \pi[2],$$

$$\pi[2]=\pi[5]=b$$

で、

$$\pi[1] = \tau[1],$$

$$\pi[2] = \tau[2] (\neq \pi[1]),$$

$$\pi[3] \neq \tau[3] (\because \tau[3] \neq \pi[1])$$

であるから、3文字以下のパターン移動では必ず照合に失敗してしまうからである。つまり、照合失敗時のパターンの移動量およびパターン中で次に照合すべき文字位置はパターンそのものから（テキストとは無関係に）あらかじめ決めておくことができる。紙面の都合でこの量を計算するアルゴリズムは割愛するが、同様な考え方は、後述するBoyer-Mooreのアルゴリズムでも使われている。

KMP法それ自身はあまり使われていないようであるが、テキストの長さに比例する時間で照合が可能であることを証明したこと、後に出現するAho-Corasick法やBoyer-Moore法に示唆を与えたことで高く評価されている。

### 2. 2. 3 Boyer-Mooreの方法<sup>1),2)</sup>

BoyerとMooreは、「パターンの右端から逆向きにテキストとの照合を行ったほうが左端から行うより多くの情報が得られる」ということを利用した方法を考案した。例えば図2. 2のような場合、パターン（“state”）の右端から照合を始めると、どちらも‘e’なので照合は成功する。この場合次に1つ手前の文字の照合を行う。今度は照合に失敗する。このときパターンを右に4以下分動かしても必ず失敗する

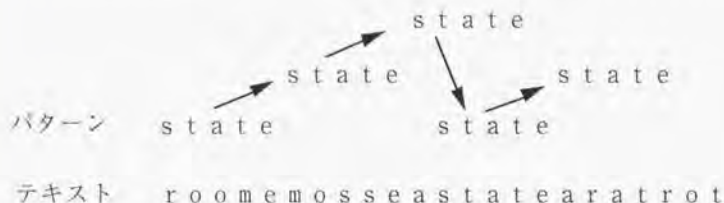


図2. 2 BM法による照合手順

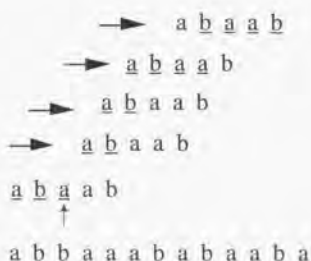


図2. 1 KMP法による照合手順

ことは、

- ・ 'm' がパターン中に存在しない
- ・ 'e' はパターンの先頭文字ではない

ことから明らかであるので、5だけパターンを右に動かしてから照合を再開してもよい、つまりパターンの右端から照合していくことによって無駄な比較を減らすことが可能である。このときパターンを動かす量はパターン自身と照合失敗時のテキスト上の文字とから決まるので、照合を始める前にあらかじめ計算しておく。これがBM法の基本的な発想である。

照合失敗時のテキスト ( $\tau$ ) の文字  $c$  から得られるパターン ( $\pi$ ) をずらす文字数  $d[c]$  は、次のように定式化できる。ただし、パターンの長さを  $m$  とする。

$$d[c] = \min\{i \mid (i=m) \text{ または } (0 \leq i < m) \text{ かつ } \pi[m-i] = c\}$$

パターン自身から得られる移動量  $dd[p]$  は、

$$dd[j] = \min\{i+m-j \mid (i \geq 1) \text{ かつ } (i \geq j \text{ または } \pi[j-i] \neq \pi[j])$$

$$\text{かつ (任意の } k (j < k \leq m) \text{ に対して } (i \geq k \text{ または } \pi[k-i] = \pi[k]))\}$$

となる。この  $d$  と  $dd$  を用いればBM法の照合アルゴリズムは図2.3のようになる。

(図2.3以下、アルゴリズムは全てPascal風に記述する。)  $d$  を求めるアルゴリズムは簡単なので省略するが(定義通りコーディングすればよい)、 $dd$  を求めるアルゴリズムは多少複雑であり、図2.4のようになる。図2.4で  $f$  は、

$$f[j] = \min\{i \mid (i=m) \text{ または } [(j < i < m) \text{ かつ}$$

$$(\pi[i+1] \cdots \pi[m] = \pi[j+1] \cdots \pi[m+j-i])]\}$$

で定義される関数であり、 $dd$  の計算のための補助として用いられている。

---

```
begin j:=m;
  repeat i:=m;
    while (i>0) and (π[i]=τ[j]) do
      begin i:=i-1; j:=j-1 end;
    if i≠0 then j:=j+max(d[τ[j]],dd[i])
  until (i=0) or (j>n)
  (i=0ならばτ[j+1]から始まる部分文字列がπ
  と一致している。i<0ならばτ中にπはない。)
end
```

---

図2.3 BM法の照合アルゴリズム



BM法のアルゴリズムの効率については、最良の場合probe rateが $1/m$ となることは明らかである。テキスト中にパターンが全く現れない場合には、最悪の場合でも4以下になることが証明されている<sup>3)</sup>。(実際にはもっと小さいであろうと予想されている。)この方法のprobe rateはランダムテキストを対象にした場合、1以下になることが証明されている<sup>1)</sup>。

---

```
begin
  for k:=1 to m do dd[k]:=2*m-k;
  j:=m; i:=m+1;
  while j>0 do
    begin
      f[j]:=i;
      while (i<m) and (π[j]≠π[i])
        begin
          dd[i]:=min{dd[i],m-j};
          i:=f[i]
        end;
      i:=i-1; j:=j-1
    end;
  t:=i; (f[0]に相当)
  for i:=1 to m do
    begin
      dd[i]:=min{dd[i],t+m-i};
      if (i<t) then t:=f[t]
    end
  end
end
```

---

図2.4 関数 $dd$ を求めるアルゴリズム

## 2. 2. 4 Aho-Corasickの方法<sup>4)</sup>

2. 2. 1から2. 2. 3までは全て照合したいパターンが1個の場合の方法であった。これらの方法で複数のパターンを照合しようとするとうパターン数だけの繰り返しが必要となってしまう。AhoとCorasickは有限オートマトンの考え方を取り入れることによって複数のパターンを同時に照合する方法を考案している。

AC法はパターンの集合からパターン照合機械(以下pmmと略記)と呼ばれる一種の有限オートマトンを作り、それをテキスト上で左から右へ1回走査させて各パターンの出現位置をことごとく検出するアルゴリズムである。pmmはgoto, failure, 出力という3つの関数から構成される。例えば、パターンが“state”, “east”, “smart”の場合のpmmは図2. 5のようになる。図で実線はgoto関数を、破線はfailure関数を、下線付きの文字列はその状態における出力を示している。破線で示した遷移は通常のオートマトンのε動作、すなわちヘッドを固定したまま状態だけを変えることを表している。ただし図2. 5では状態5, 8, 9以外から状態0への破線は省略してある。

pmm作成のアルゴリズムの記述は省略するが、後述するFAST法でも同様なアルゴリズムが用いられる。このpmmをテキスト上で走査させれば、その状態遷移によってパターンの存否がわかることになる。テキストが“osseastatear”の場合を例に状態遷移の様子を示すと図2. 6のようになる。この例からも分かるようにAC法ではパターンの数に関わりなくテキストの長さに比例する時間で照合が

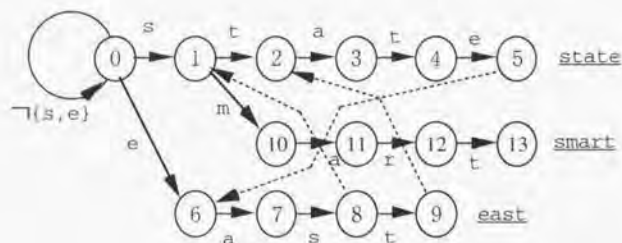


図2. 5 AC法におけるパターン照合機械

可能である。failure関数をなくし、goto関数だけのpmmを作ることも可能で、この場合にはprobe rateを完全に1にすることができる。

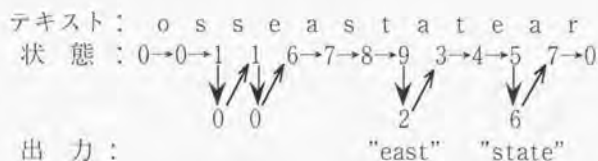


図 2. 6 pmmの状態遷移

## 2. 3 FAST法のアルゴリズム

前述したように、パターンが1つの場合の照合はBoyer-Moore法が最も効率の良い方法として知られている。ところが、複数のパターンの照合にはパターンの数だけの繰り返し使用が必要となって、急激に効率が低下してしまう。一方、Aho-Corasick法は複数パターンを同時に扱うことが可能である。AC法は、KMP法に有限オートマトンの考えを取り入れて複数パターンを同時に照合できるようにしたものを見なすことができる。BM法に有限オートマトンの考えを取り入れれば、複数パターンの場合でも平均的に亜線形 (sublinear, probe rate が1以下) の照合が可能となるはずである。筆者はBM法とAC法の両者の考え方を取入れることによって、複数パターンの場合でも効率的にパターンを照合できるアルゴリズムを考案した。この方法をFAST法 (a Flying Algorithm for Searching Terms) と呼ぶことにする。以下でFAST法の考え方とアルゴリズムについて述べる。

### 2. 3. 1 FAST法の考え方

BM法の「右端からパターンを照合する」というアイデアとAC法の「pmmによって複数パターンの同時照合を可能にする」というアイデアを両方とも組み込むことができれば、複数パターンにおいてもBM法と同様な低いprobe rateを実現できるはずである。

パターンの集合は図2.5の場合と同じで、テキストは図2.2と同じであると仮定する。このとき複数のパターンの照合を同時に行うにはパターンの右端を揃え

て図2. 7のようにテキストの上に並べてから照合を開始すればよい。図2. 7の場合、テキスト上の文字 'm' は3つのパターンの右端のどれとも一致しないからパターン群を後ろに動かすことができる。'm' はパターン "smart" の中にしか現れないから3文字分動かしても探索もれが起きる心配はない。したがってパターン群を3文字分右へ動かしてから照合を再開すればよいことになる。問題は、探索もれが生じない範囲でパターン群を動かせる最大の量をどのようにすれば求められるか、ということにある。

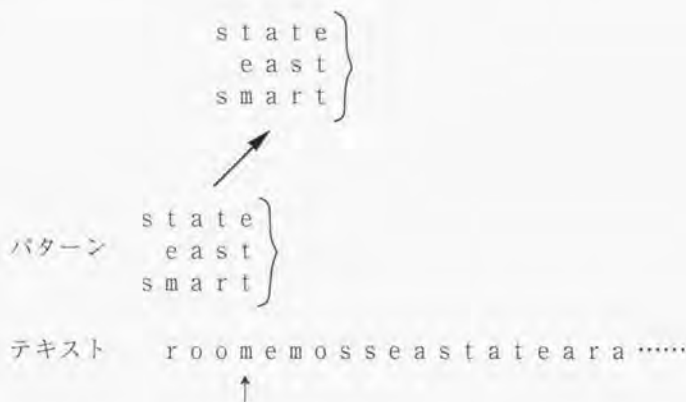


図2. 7 複数パターンの右端からの照合

これには図2. 8に示す2つの場合を考慮すればよい。つまりテキスト上の文字  $a_i$  から始まった照合があるパターン  $\alpha$  と途中まで成功して、 $a_i$  で失敗した場合、

・テキストの部分列  $a_i a_{i+1} \dots a_j$  があるパターン  $\beta$  の部分列になっている場合 (a) と

・テキストの部分列  $a_i a_{i+1} \dots a_j$  ( $i < k \leq j$ ) があるパターン  $\beta$  の先頭の部分列になっている場合 (b)

である。この2つの場合、パターン  $\beta$  の一致する部分列をテキストの当該部分列に重ねるような位置にパターン群を動かして (図2. 8のシフト量  $S_{\alpha\beta}$ ) 照合操作を再開してもパターン  $\beta$  の探索もれが起こる可能性はない。すなわち、 $S_{\alpha\beta}$  を厳密に

定義すれば、(ただし、パターン  $a$  の長さを  $m$  とし、パターン  $\beta$  の長さを  $m'$  とする。)

$$S_{a\beta} = \min \{ d \mid [ (d \geq 1) \text{ かつ } (m' - d) \leq (j - i) \text{ または } \beta_{m' - d - j + i} = a_i) \text{ かつ (任意の } n (m - j + i < n \leq m) \text{ に対して}$$

$$( (m' - d) \leq (m - n) \text{ または } \beta_{m' - d - m + n} = a_n ) ] \text{ または } [ d = m' ] \}$$

となる。

もちろん、場合 (a)、(b) に相当する部分列がない場合にはパターン  $\beta$  の長さだけ動かすことが可能である。このシフト量  $S_{a\beta}$  はパターンの集合が与えられた時点で決まってしまう。

BM法でシフト量を求める方法も基本的にはこの考え方と同じである。ただBM法ではパターン  $a$  とパターン  $\beta$  は同一であるのに対し、複数パターンを考慮するためには、照合中のパターン  $a$  に対し、すべてのパターン (同一パターンを含む) に

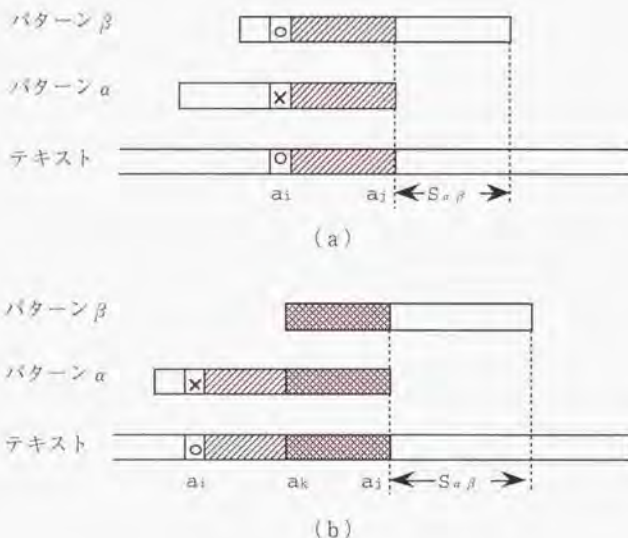


図 2. 8 パターン群シフト量の計算

対する  $S_{\alpha, \beta}$  を求め、その最小値  $S_m$  つまり、

$$S_m = \min |S_{\alpha, \beta} : \beta \in \text{パターン集合}|$$

をパターン群のシフト量としなければならない。また、BM法では照合に失敗した時、テキスト上の文字 ( $a_j$ ) に対するシフト量を  $d$  という関数で、それまでの照合情報 ( $a_{1,1}, a_{1,2}, \dots, a_{1,j}$ ) に基づくシフト量を  $dd$  という関数で別々に取り扱っていたが、FAST法では  $a_{1,1}, a_{1,2}, \dots, a_{1,j}$  の出現をまとめて取り扱っている点でも異なる。このため、パターンが1つの場合でもFAST法の方がBM法よりシフト量が多くなる場合が生じ、一般に効率が良い。

一方、照合に成功している間は、どのパターンのどの部分を照合しているかの情報はAC法と同様に  $pmm$  の状態として保持しておけばよい。ただAC法と異なる点は状態遷移を引き起こす文字入力テキスト上で右から左への方で取られるだけである。この状態遷移を求めるアルゴリズムは2.3.2で述べるが、状態遷移のための  $pmm$  は前述のシフト量の計算にも使われる。

ここで述べるアルゴリズムは、簡単に言えば、照合に成功している間はAC法と同様の状態遷移を行い、照合失敗時の動作には複数パターンに拡張したBM法のシフト関数を利用するというものである。

### 2.3.2 FAST法のアルゴリズム

図2.9にgoto関数と出力関数を求めるためのアルゴリズムを示す。このアルゴリズムは  $pmm$  をパターンの後ろから作る点を除けばAC法と同じである。ただし、図2.10のアルゴリズムで用いるために、節点0からの深さを示す関数  $depth$  と、パターンの左端に対応する状態を示す関数  $emap$  と、パターンの最小長  $minlen$  も計算している。

-----

パターンの集合を  $\{y_1, y_2, \dots, y_n\}$  とする。

状態  $s$  が作られるとき、 $out(s) = \text{empty}$  と仮定する。

$g(s, a)$  の初期値は **fail** と仮定する。

出力はgoto関数  $g$  と出力関数  $out$  である。

図2.10のアルゴリズムで利用する  $minlen, depth, emap$  もここで計算される。

```

begin
  nst:=1; depth(0):=0;
  minlen:=large; /* large はパターン長に比べて充分大きい正整数 */
  for i:=1 to k do enter( yi, i)
end

procedure enter(a1a2...am, k)
begin
  state:=0; j:=m;
  while ( g (state, aj ) ≠ fail) and ( j > 0 ) do
    begin
      state:=g (state, aj);
      j:=j-1
    end
  while j > 0 do
    begin
      g (state, aj):=nst;
      depth(nst):=m-j+1;
      state:=nst;
      nst:=nst+1;
      j:=j-1
    end
  emap(k):=state;
  out(state):={a1a2...am};
  minlen:=min( minlen, m)
end

```

---

図 2. 9 goto 関数を求めるアルゴリズム

図2.10には図2.9のアルゴリズムの結果を利用してfailure関数を求めるアルゴリズムを示す。図2.10で $\phi$ はAC法のfailure関数に相当するものであり、状態 $st$ の場合、 $depth(\phi(st))$ はBM法のアルゴリズムでシフト量( $sd$ )を求める際に用いた補助関数 $f$ に相当している。 $\phi$ を利用してFAST法のfailure関数 $F$ を計算する。図2.10のステートメントAで $F$ を更新しているのは図2.8の(a)の場合に相当している。つまり、パターンを照合していて、状態 $fst$ -入力 $c$ で照合に失敗した場合、次のポイントをパターン $\beta$ の最後尾に移すための計算である。関数 $dcm$ は図2.8の(b)に相当するシフト量を表わしている。 $dcm$ の計算の最初のループで、照合済みのパターン $a$ のsuffix(後綴り)をprefix(前綴り)に持つパターン $\beta$ を探して、その最後尾までのポイントのシフト量(の最小値)を求める。後のループでは(未決定のものを含めて) $dcm$ の値をその1つ手前の状態(節点)から順に決めていくことによって、 $dcm$ に矛盾(後ろの状態の $dcm$ が直前の状態の $dcm$ より2つ以上大きくなること)がないようにしている。アルゴリズムの最後で $g(j,c)=fail$ となる $F(j,c)$ に対して、 $\min(F(j,c), dcm(j))$ をセットして、failure関数 $F$ が完成する。

---

goto関数 $g$ と $minlen, depth, emap$ は図2.9のアルゴリズムによって求められているとする。

$F(s, a)$ と $dcm(s)$ の初期値は**large**(十分大きな数)とする。  
出力はfailure関数 $F$ と出力関数 $out$ である。

```

begin
  queue:=empty;   $\phi(0):=-1$ ;
  for each c such that  $g(0,c)=s \neq fail$  do
    begin
      queue:=queue.s;
       $\phi(s):=0$ 
    end;

```



```

while queue#empty do
  begin
    let queue:=r.tail;
    queue:=tail;
    for each c such that g(r,c)=s#fail do
      begin
        queue:=queue.s;
        fst:=  $\phi$ (r);
        while (fst>=0) and (g(fst,c)=fail) do
          begin
            f(fst,c):=min{ f(fst,c), depth(r)}; <--- A
            fst:=  $\phi$ (fst);
          end;
          if fst>=0 then  $\phi$ (s):=g(fst,c);
            else  $\phi$ (s):= 0
          out(s):=out(s).out( $\phi$ (s))
        end
      end;
    end;
  for j:=1 to k do
    begin
      fst:= $\phi$ (emap(j));
      while fst > 0 do
        begin
          ddm(fst):=min{ ddm(fst), depth(emap(j))};
          fst:= $\phi$ (fst)
        end
      end;
    end;
  ddm(0):=minlen;
  queue:={0};

```

```

while queue ≠ empty do
  begin
    let queue := r.tail;
    queue := tail;
    for each c such that g(r,c) = s ≠ fail do
      begin
        queue := queue.s;
        ddm(s) := min{ ddm(s), ddm(r)+1 };
      end
    end;
  end;
for j := 0 to nst-1 do
  for each c such that g(j,c) = fail do
    f(j,c) := min{ f(j,c), ddm(j) };
  end;
end

```

図2. 10 failure関数を求めるアルゴリズム

図2. 11に、パターンが図2. 5と同じ場合、図2. 9、図2. 10のアルゴリズムで求められるFAST法のpmmを示す。実線、破線、下線付き文字列の意味は図2. 5と同じである。ただし、破線は関数 $f$ でなく関数 $\phi$ に相当している。

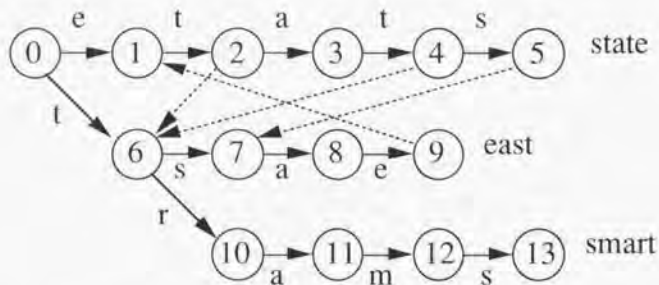


図2.11 FAST法におけるパターン照合機械

ポインタの移動の向きが異なる点を除けば、FAST法の照合に成功している間の動作はAC法と全く同じである。しかし、照合失敗時の動作は異なり、FAST法ではfailure関数によって、ポインタのシフト量を求める。ところで、failure関数はgoto関数がfailとなる場合にしか定義されないからfailure関数の値を負に変えることによってgoto関数と同一の関数にまとめることができる。この様にして求めたgoto関数 $g$ を表2.1に示す。

表2.1 goto関数の例

文字種/状態	0	1	2	3	4	5	6	7	8	9	10	11	12	13
a	-2	-4	3	-6	-7	-8	-2	8	-6	-7	11	-7	-8	-9
e	1	-4	-5	-6	-7	-8	-5	-5	9	-7	-6	-7	-8	-9
m	-3	-4	-5	-6	-7	-8	-5	-5	-6	-7	-6	12	-8	-9
r	-1	-4	-5	-6	-7	-8	10	-5	-6	-7	-6	-7	-8	-9
s	-1	-4	-5	-6	5	-8	7	-5	-6	-7	-6	-7	13	-9
t	6	2	-5	4	-7	-8	-5	-5	-6	-7	-6	-7	-8	-9
(other)	-4	-4	-5	-6	-7	-8	-5	-5	-6	-7	-6	-7	-8	-9

こうして求めたgoto関数 $g$ を用いたパターン照合のアルゴリズムは図2.12のようになる。すなわち、AC法ではテキスト上のポインタは常に前方へ1だけ増加するが、FAST法では照合に成功している間は後方へ1だけ戻り、失敗したときは $g$ の絶対値分だけ前方に進む。

---

テキストの文字列 $x$ を $x=a_1a_2\cdots a_n$ とする。  
goto関数を $g$ 、出力関数を $out$ とする。  
出力は $x$ 中に出現するパターンとその位置である。

```
begin
  aq:=umc; /* umc は決して照合に成功しない文字 */
  g:=minlen; /* minlen はパターンの最小文字長 */
  while q ≤ n do
    begin
      state:=0;
      while g(state, aq) > 0 do
        begin
          state:=g(state, aq);
          if out(state) ≠ empty then print q, out(state);
          q:=q-1;
        end
      q:=q-g(state, aq);
    end
  end
```

---

図 2. 1 2 文字列照合アルゴリズム

図 2. 1 3 に図 2. 2 と同じテキストを用い、FAST法の照合動作の例を示す。図 2. 1 3 でstate は pmmの状態を、orderは照合の順番を示している。状態0のときテキスト上の文字が't'であれば状態6に遷移し、探索対象から"state"が除外され、次に1つ手前のテキスト上の文字が照合される。照合に失敗すると状態とテキスト上の文字によって決まる量だけパターン群が右に移動する。こうして四角で囲んだところでパターンが検出される。図 2. 1 3 からわかるようにテキス

ト上の文字すべてを照合しないですむ場合がある。反対にテキスト上の同一文字が何回も照合されることもある。例えば図2.13のテキスト“...seast...”の‘a’は3回照合される。

```

state:           5 4 3 2 1 0
order:           1413121110 9

state:           9 8 7 6 0
order:           8 7 6 5 4

state:           0   0   0
order:           1   2   3

テキスト   r o o m e m o s s e a s t a t e a r a .....

```

図2.13 FAST法による照合手順

## 2.4 あとがき

全走査法に属するいろいろな全文検索手法について解説し、BM法の「右端からの照合」というアイデアとAC法の「パターン照合機械」のアイデアを結合することによって、複数パターンを同時に高速に探索する方法（FAST法）を示した。

実は、FAST法はパターン数が1つの場合にもBM法より効率が良い。その理由は、テキストの $a_j$ から後向きに $a_1$ まで照合して行って、そこで失敗した場合、BM法では $a_{i-1}a_{i-2}\cdots a_j$ と一致する可能性がある位置までパターンを移動するわけだが、FAST法では $a_2a_{1+1}\cdots a_j$ と1文字分長くどっているためシフト量を大きくすることができるからである。FAST法がBM法より劣っている点はFAST法の方がパターンのシフト量を計算するために保持しておくための関数の記憶容量が大きくなることである。BM法では関数 $d$ 、関数 $da$ の為に「パターン長+文字種の大きさ」の記憶容量を用意すればよいのに対し、FAST法のgoto関数は「パターン長×文字

種の大きさ」の記憶容量を必要とする。このことは、文字種が多い場合には問題になりうる。記憶容量の削減法については第3章でふれる。

F A S T法がBM法より実際の実行時間においても速い理由がもう一つある。それは、BM法では照合に失敗するたびに関数 $d$ と関数 $dd$ の大きい方を選ぶ演算を必要とするのに対し、F A S T法では1回の表探索で済むからである。そこで、文字種が多い場合にはより効果が期待できる関数 $d$ だけを用いることによって、BM法の実際の実行時間を短縮する方法(BMH法)がHorspoolによって提案されている<sup>5)</sup>。

複数パターン的高速探索法についてはKowalskiとMeltzerも提案している<sup>6)</sup>。この方法は単純に言えば、パターン長が不揃いの場合、まず最短のパターンに合わせて後を切り揃え、それらを探索する。そして、切り揃えられた各パターンそれぞれにBM法の関数 $d$ 、関数 $dd$ を求めておき、照合に失敗したときはそれらを用いてパターン群をシフトする。この方法はかなり煩雑であり、そのうえ明らかにF A S T法よりも効率が悪い。

Commentz-Walterも別の複数パターン探索法を提案している<sup>7)</sup>。この方法は、シフト量の計算に $p$  mmを使うなど、F A S T法とかなり似ているが、関数 $d$ と関数 $dd$ を別々に求めていて、「 $a_{1+1}a_{1+2}\dots a_j$ の情報しか用いない」こと、「照合に失敗するたびに関数 $d$ と関数 $dd$ の大きい方を選ぶ演算を必要とする」ことなどBM法と同じ欠点を有しているので、明らかにF A S T法よりも効率が悪い。そのうえ、F A S T法では最大の問題点は、探索漏れが生じない範囲でパターン群を動かせる最大の量をどのようにすれば求められるか、という点にあったが、Commentz-Walterのアルゴリズムにはこの計算法に間違いがあり、探索漏れが生じる可能性がある<sup>8)</sup>。

F A S T法の効率の良さはアルゴリズムそのものから直感的には明らかであろうが、第3章で、効率について理論的に考察する。また、アルゴリズムの正当性を証明し、実験によって実際に高い探索効率が得られることを示す。

## 第3章 FAST法の性質

### 3.1 まえがき

第2章では、いろいろな全文検索手法について解説し、複数パターンを同時に高速に探索する方法（FAST法）の考え方と具体的なアルゴリズムについて示した。そして、パターン数が1つの場合でさえもFAST法の方がBM法より効率が良いことも述べた。FAST法の効率の良さはアルゴリズムそのものから直感的には明らかであるが、この章で効率について理論的に考察する。

いくら高速な探索方法であっても、常に正しく動作するものでなければ価値はない。この章ではFAST法のアルゴリズムの正当性についても証明する。さらに実験によって実際に高い探索効率が得られることを示す。

FAST法がBM法より劣っている点はgoto関数の記憶容量が大きくなることである。この削減方法についても検討を加える。

### 3.2 FAST法の効率<sup>1)</sup>

アルゴリズムの効率を考える場合には、もちろん主動作の効率だけでなく、前処理に要する時間も考慮しなければならない。2.2.1で述べた「素朴な方法」は全く前処理時間を必要としない。KMP法やAC法では $m$ を（全）パターンの長さとする $O(m)$ 、BM法では $O(q+m)$ （ $q$ は文字種の数）の前処理時間を必要とする。これに対してFAST法では $O(mq)$ の前処理時間を必要し、他の方法と比べて劣っている。しかし、通常対象とするファイルはかなり大きいと考えられるので、（そうでなければ「素朴な方法」で十分である）前処理時間は主動作に要する時間に比べて無視できると仮定する。したがって、以下では文字列照合操作に直接関わる時間のみを考える。

文字列照合アルゴリズムの優劣は一般にprobe rate（テキスト1文字当りの照合回数）で評価される。この場合は、一般的にいう効率はprobe rateの逆数と考えてよい。

FAST法のprobe rateを $Q$ とすると、明らかに、最良時には

$$Q = 1 / \text{minlen} \quad (\text{minlen: パターンの最小長})$$

となる。最悪時としては照合の成功する確率が最大の場合を想定すれば良い。この

場合でもパターンが1つ見つかり、その直後で失敗する毎に最低1つは照合開始位置が右にずれるので、Qは

$$Q = \text{maxlen} + 1 \quad (\text{maxlen: パターンの最大長})$$

となる。すなわち最悪時でもQがパターン数に依存しないという望ましい特徴を有している。

テキストもパターンもq種の文字集合からランダムに抽出されていると仮定して、Qの期待値を求めてみる。ただし、パターンの長さはすべてmとする。

今、パターンが1つの場合を考えて

テキストを  $\tau_1 \tau_2 \cdots \tau_j \cdots \tau_n$

パターンを  $\pi_1 \pi_2 \cdots \pi_s \cdots \pi_m$

で表すものとする。jがmより十分大きくて $\tau_j$ から始まるパターン照合の成功する確率がjに依存しない状態(定常状態)を考える。 $P_{s,i}$ をs回の照合成功の後s+1回目で失敗して最初の照合位置からiだけずれたところから照合が再開される確率とすると、 $\tau_j$ と $\pi_m$ から始まった照合は $\tau_{j+s+1}$ と $\pi_{m-s+1}$ との照合までは成功し、 $\tau_{j+s}$ と $\pi_{m-s}$ の照合で失敗する。この時、パターンのシフト量がiとなることから、 $\tau_{j+s} \tau_{j+s+1} \cdots \tau_j$ は $\pi_{m-i}$ から(後向きに)始まる照合では成功するが、それ以前から始まる照合で成功することはない。したがって、

s=0のとき

$$P_{0,i} = P\left(\left(\tau_j \neq \pi_m\right) \wedge \left(\tau_j \neq \pi_{m-1}\right) \wedge \cdots \wedge \left(\tau_j \neq \pi_{m-i+1}\right) \wedge \left(\tau_j = \pi_{m-i}\right)\right)$$

0 < s のとき

$$P_{s,i} = P\left(\left(\tau_j = \pi_m\right) \wedge \left(\tau_{j-1} = \pi_{m-1}\right) \wedge \cdots \wedge \left(\tau_{j-s+1} = \pi_{m-s+1}\right) \wedge \left(\tau_{j-s} \neq \pi_{m-s}\right)\right)$$

$$\wedge \left(\tau_{j-s} \tau_{j-s+1} \cdots \tau_j \neq \pi_{m-s-1} \pi_{m-s} \cdots \pi_{m-1}\right)$$

$$\wedge \left(\tau_{j-s} \tau_{j-s+1} \cdots \tau_j \neq \pi_{m-s-2} \pi_{m-s-1} \cdots \pi_{m-2}\right)$$

⋮

$$\wedge \left(\tau_{j-s} \tau_{j-s+1} \cdots \tau_j \neq \pi_{m-s-i+1} \pi_{m-s-i} \cdots \pi_{m-i+1}\right)$$

$$\wedge \left(\tau_{j-s} \tau_{j-s+1} \cdots \tau_j = \pi_{m-s-i} \pi_{m-s-i+1} \cdots \pi_{m-i}\right) \quad (1)$$

ここで、 $\tau_{j-s} \tau_{j-s+1} \cdots \tau_j \neq \pi_{m-s-1} \pi_{m-s} \cdots \pi_{m-1}$ 等は

$(\tau_{j-s} \neq \pi_{m-s-1}) \vee (\tau_{j-s+1} \neq \pi_{m-s}) \vee \cdots \vee (\tau_j \neq \pi_{m-1})$ 等の略記であり、

$\tau_{j-s} \tau_{j-s+1} \cdots \tau_j = \pi_{m-s-i} \pi_{m-s-i+1} \cdots \pi_{m-i}$ は



$(\tau_{j-s} = \pi_{m-s-i}) \wedge (\tau_{j-s+1} = \pi_{m-s-i+1}) \wedge \cdots \wedge (\tau_j = \pi_{m-i})$  の略記である.

( $0 < s < (m-1)$  で  $i < (m-s)$  の場合)

ただし,  $(m-s) \leq i$  のときは  $m-s-i \leq 0$  となるから (1) 式中の  $\pi_{m-s-i}$  等に対応するものがなくなるので, その分を評価しなくてよいため P 中の後半の項は短くなる.

$s = m-1$  や  $s = m$  のときも同様で, 対応するものがない分, 上述の式は短くなる.

ところで, (1) 式の P 中の各項は独立でない. 例えば  $m=3, s=1, i=2$  のときは,  $\tau_j = \tau_{j-1}$  の場合と,  $\tau_j \neq \tau_{j-1}$  の場合を分けて考えれば

$$\begin{aligned} & P((\tau_j = \pi_3) \wedge (\tau_{j-1} \neq \pi_2) \wedge (\tau_{j-1}\tau_j \neq \pi_1\pi_2) \wedge (\tau_j = \pi_1)) \\ &= P((\tau_j = \tau_{j-1}) \wedge (\tau_j = \pi_3) \wedge (\tau_{j-1} \neq \pi_2) \wedge (\tau_{j-1}\tau_j \neq \pi_1\pi_2) \wedge (\tau_j = \pi_1)) \\ &+ P((\tau_j \neq \tau_{j-1}) \wedge (\tau_j = \pi_3) \wedge (\tau_{j-1} \neq \pi_2) \wedge (\tau_{j-1}\tau_j \neq \pi_1\pi_2) \wedge (\tau_j = \pi_1)) \\ & \quad [(\tau_j = \tau_{j-1}) \wedge (\tau_{j-1} \neq \pi_2) \text{ のときも, } (\tau_j \neq \tau_{j-1}) \wedge (\tau_j = \pi_1) \text{ のときも} \\ & \quad \tau_{j-1}\tau_j \neq \pi_1\pi_2 \text{ となるから}] \\ &= P((\tau_j = \pi_3) \wedge (\tau_{j-1} \neq \pi_2) \wedge (\tau_j = \pi_1)) \\ &= (q-1)/q^3 \end{aligned}$$

となつて, 各項の確率の積  $((q-1)(q^2-1)/q^5)$  とはならない. しかし,  $q \gg 1$  の場合はこの例からもわかるように各項を独立と考えてもその誤差は小さい<sup>\*</sup>.

☆ 簡単に述べると,  $s=0$  の場合は各項は独立となるため誤差は生じない.

$s > 0$  のときは,  $(\tau_{j-s}\tau_{j-s+1}\cdots\tau_j \neq \pi_{m-s-k}\pi_{m-s-k+1}\cdots\pi_{m-k})$  のすべてを

$(\tau_j \neq \pi_{m-k})$  で置き換えれば P の下限 ( $P_d$ ) が, 1 で置き換えれば P の上限 ( $P_u$ )

が得られる. 各項を独立と見なした P の推定値 ( $\hat{p}$ ) は  $P_d \leq \hat{p} \leq P_u$  となり,

絶対誤差は

$$\begin{aligned} (m-s) \leq i \text{ のとき} & \quad \frac{i-1}{q^{s+m-i+1}} \\ (m-s) > i \text{ のとき} & \quad \frac{i-1}{q^{2s+2}} \end{aligned}$$

より小さくなるので, 「各項を独立と考えてもその誤差は小さい」としてよい.

もう 1 つ考慮すべきことがある. 例えば 1 つ前が  $P_{0,1}$  から  $\tau_j$  の照合に入る場合は  $\tau_{j-1} = \pi_2$  となるから  $P_{1,2}$  等で P 中の項を満たさない.  $\tau_j$  から始まる照合は一般に  $(m-1)$  前の状態まで考えなくてはならないが, 1 つ手前の状態だけを考慮に入れる

こととする。つまり、時点  $k$  において状態  $(x_k)$  が  $\psi_j$  となる確率は、(最悪の場合) 時点  $k-1, k-2, \dots, k-m+1$  の  $(m-1)$  個の状態に依存するので

$$P(x_k = \psi_j | x_{k,m+1} = \psi_{i,k,m+1}, x_{k,m+2} = \psi_{i,k,m+2}, \dots, x_{k-1} = \psi_{i,k-1})$$

と  $(m-1)$  重マルコフ連鎖を考えなければならないが、これを  $P(x_k = \psi_j | x_{k-1} = \psi_i)$  で近似するわけである。(2つ以上前の状態の影響は小さいのでこうしても誤差は少ない。なぜなら、 $i=m$  のときはテキスト上で同じ場所の文字が照合されないから以前の状態に影響されないし、 $i \neq m$  のときは、 $P_{s,i} < 1/q$  だから2つ以上前の影響が  $1/q^2$  を超えることはない。)

前述のように単純化のために  $P$  中の各項を独立であると仮定しても誤差は小さいので、(1) 式は

$$P_{s,i} = (t_s - t_{s+1}) \left( \frac{q^{s+1} - 1}{q^{s+1}} \right)^{i-1} \frac{1}{q^{s+1}}$$

と近似できる。

ここで  $t_s$  は  $\tau_j$  から始まる照合が  $s$  回照合に成功する確率であり、

$$t_s = P((\tau_j = \pi_m) \wedge (\tau_{j-1} = \pi_{m-1}) \wedge \dots \wedge (\tau_{j-s+1} = \pi_{m-s+1}))$$

である。上述のようにこの  $P$  は以前の状態に影響され各項の単純の積とはならない。 $\tau_j$  から始まる照合の1つ手前の状態だけを考慮すれば、例えば  $P_{2,2}$  から  $\tau_j$  の照合に入る場合には(パターン長が5以上の場合)5回照合に成功する確率は  $\tau_{j,2}, \tau_{j,3}, \tau_{j,4}$  の照合成功は予め確定しているから、 $p_5$  ではなく条件付き確率の  $p_5/p_3$  となる。したがって、 $t_s$  (の近似値) は、

$$\begin{aligned} t_s = & \left( 1 - \sum_{h=0}^m P_{h,1} - \dots - \sum_{h=0}^m P_{h,s-1} \right) p_s \\ & + \left( \sum_{h=0}^m P_{h,s-1} + \sum_{h=1}^{s-2} P_{0,h} \right) p_s / p_{s-1} \\ & \vdots \\ & + \left( \sum_{h=k-1}^m P_{h,s-k} + \sum_{h=1}^{s-k-1} P_{k-1,h} \right) p_s / p_k \\ & \vdots \\ & + \left( \sum_{h=s-2}^m P_{h,1} \right) p_s / p_{s-1} \end{aligned}$$

となる。ここで、

$$p_s = 1/q^s$$

であり、 $p_s/p_k$  は  $k$  回の照合成功が分っている場合  $s$  回照合が成功する条件付き確率である。また、 $t_0 = 1$ ,  $t_j = p_j$ ,  $t_{m+1} = 0$  である。

パターンが  $k$  個のときはパターン中のある連続した  $s$  個の文字がテキストの現在照合している部分と一致する確率  $p_{s,i}$  は

$$p_s = 1 - (1 - 1/q^s)^k$$

となるので、以上をまとめると、(単純化のためパターン長はすべて  $m$  とする)

$P_{s,i}$  の一般式は

●  $0 \leq s \leq (m-1)$  のとき

$i < (m-s)$  の場合

$$P_{s,i} = (t_s - t_{s+1}) p_{s+1} (1 - p_{s+1})^{i-1}$$

$i \geq (m-s)$  の場合

$$P_{s,i} = (t_s - t_{s+1}) p_{m-i} (1 - p_{s+1})^{m-s-1} \prod_{k=m-i+1}^s (1 - p_k)$$

●  $s \leq m$  のとき

$$P_{s,i} = (t_s - t_{s+1}) p_{m-1} \prod_{k=m-i+1}^{m-1} (1 - p_k)$$

(ただし  $\prod_{k=a}^b x_k$  は  $a > b$  のときは 1 とする)

となる。

これらの  $P_{s,i}$  ( $0 \leq s \leq m$ ,  $1 \leq i \leq m$ ) と  $t_j$  ( $1 \leq j \leq m$ ) に関する連立方程式を解けば  $P_{s,i}$  と  $t_j$  のそれぞれの値が求まる。

位置  $\tau_j$  から始まる照合は失敗するまで延べ  $\sum_{j=0}^m t_j$  回の照合を必要とするが、失敗

したときには  $\sum_{i=1}^m i \left( \sum_{s=0}^m P_{s,i} \right)$  だけ次の照合開始位置が進むので probe rate  $Q$  は

$$Q = \sum_{j=0}^m t_j / \sum_{i=1}^m i \left( \sum_{s=0}^m P_{s,i} \right) \quad (2)$$

となる。

$q$  が  $m, k$  に対して十分大きい ( $q \gg mk$ ) ときには  $r = 1/q$  とすると

$$p_j = kr^j$$

となるから

$$P_{0,i} = (t_0 - t_1)kr \quad (1 \leq i \leq (m-1) \text{ のとき})$$

$$P_{0,m} = (t_0 - t_1)(1 - (m-1)kr)$$

$1 \leq s \leq (m-1)$  のとき

$$P_{s,i} = (t_s - t_{s+1})kr^{s+1} \quad (1 \leq i \leq (m-s-1) \text{ のとき})$$

$$P_{s,i} = (t_s - t_{s+1})kr^{m-i} \quad ((m-s) \leq i \leq (m-1) \text{ のとき})$$

$s = m$  のとき

$$P_{m,i} = t_m kr^{m-i} \quad (1 \leq i \leq (m-1) \text{ のとき})$$

$$P_{m,m} = t_m(1 - kr)$$

$$t_j = kr^j \quad (1 \leq j \leq m) \text{ となるから}$$

$$\sum_{j=1}^m t_j = 1 + kr + kr^2 + \dots + kr^m = 1 + kr$$

$$\sum_{i=1}^m i \left( \sum_{s=0}^m P_{s,i} \right) = m - \frac{m^2 - m}{2} kr$$

したがって

$$Q \approx \frac{1 + kr}{m \left( 1 - \frac{m-1}{2} kr \right)} \approx \frac{(1 + kr) \left( 1 + \frac{m-1}{2} kr \right)}{m}$$

$$\left( q \gg mk \text{ より } \frac{m-1}{2} kr \ll 1 \text{ だから} \right)$$

$$= \frac{1}{m} + \frac{(m+1)k}{2mq} \quad (3)$$

となる。

### 3. 3 正当性の証明<sup>2),3)</sup>

アルゴリズムの正当性を証明するには、健全性、完全性、停止性を証明すればよい。文字列探索アルゴリズムに即して言えば、これらは

健全性：間違ったパターンを出力しないこと

完全性：対象テキスト中に含まれるパターンを漏れなく出力すること

停止性：無限ループに入ったりすることがなく、必ず停止すること

となる。

まず、停止性は図2. 10のアルゴリズムにより、どの  $f(j, c)$  も  $depth(j)$  より大きくなることから明らかである。健全性についても、図2. 9のアルゴリズムに基づく状態遷移は図2. 11のオートマトン上を動くだけであるから、関係ないパターンを検出することは決してないことがわかる。

完全性の証明のためには準備が必要である。証明に必要な集合や写像、関係などの定義や関数の特徴付けを行ないつつ、完全性の証明を与える。

### 3. 3. 1 完全性の証明

文字列  $w$  の prefix を  $PRE(w)$ 、suffix を  $SUF(w)$  とする。つまり、アルファベット集合を  $\Sigma$  としたとき、 $\Sigma^*$  の任意の要素  $w$  に対し、

$$PRE(w) = \{x \in \Sigma^* \mid \exists y \in \Sigma^* : xy = w\}$$

$$SUF(w) = \{y \in \Sigma^* \mid \exists x \in \Sigma^* : xy = w\}$$

と定義する。文字列集合  $X$  に対しては、

$$PRE(X) = \bigcup_{w \in X} PRE(w), \quad SUF(X) = \bigcup_{w \in X} SUF(w)$$

とする。パターンの集合を  $K = \{a_1, a_2, \dots, a_k\}$ 、状態遷移グラフに現われる状態全体の集合を  $S$  とし、写像  $state : SUF(K) \rightarrow S$  を次式で定義する。

$$\begin{cases} state(\varepsilon) = 0, \\ state(\sigma w) = g(state(w), \sigma) \quad (w, \sigma w \in SUF(K), \sigma \in \Sigma) \end{cases}$$

$S$  上の二項関係  $\vdash_\phi$  を次のように定義する。 $s, t$  を  $S$  の任意の要素とすると、

$$s \vdash_\phi t \iff \phi(s) = t$$

とする。

図2. 9で求めた関数  $depth, \text{emap}$  と上記の定義には以下の関係がある。

ただし、 $u$  は  $S$  の任意の要素で、 $st = state(u)$  である。

$$depth(st) = |u|$$

$$\text{emap}(j) = state(a_j)$$

さて、完全性を証明するには次にあげる定理1を用いて、その後の補題1を証明すればよい。

定理1.  $u \in SUF(K)$ 、 $a \in \Sigma$ 、 $au \notin K$  とする。いま、

- (1)  $auv \in SUF(K)$
- (2)  $\exists y \in SUF(u) : yv \in K$

のいずれかを満たす  $v \in \Sigma^+$  のうち長さが最小となる  $v$  を選び、 $st = state(u)$  とおけば、

$$f(st, a) = |u| + |v|$$

である。

この定理の証明を与えるためには、さらに3つの補題を証明しなければならないが、それを後回しにして、先に補題1を証明する。

補題1. テキストの文字列を  $a_1 a_2 \dots a_n$  とし、 $a_0 = \mathbf{umc}$  (決してパターンに含まれない文字) とする。いま、

$$u = a_{p+1} \dots a_q \in \text{SUF}(K) \quad (1 \leq p \leq q \leq n)$$

$$st = state(u)$$

とおく (ただし、 $p=q$  のときは  $u = \epsilon$  とする)。

また、 $g(st, a_p) = \text{fail}$  (または0) であるとする。このとき、次の式が成り立つ。

$$[a_1 \dots a_j \in K (1 \leq i \leq j \leq n), q < j] \Rightarrow p + f(st, a_p) \leq j$$

(これは、探索漏れが起きないことを示している)

証明

(1)  $i \leq q$  のとき、 $v = a_{q+1} \dots a_j$  とおく。定理1より、

$$\begin{aligned} f(st, a_p) &\leq |u| + |v| \\ &= (q-p) + (j-q) \\ &= j-p \end{aligned}$$

である。

(2)  $q < i$  のとき、 $v = a_1 \dots a_j$  とおく。 $v \in K$  であるから、定理1より、

$$\begin{aligned} f(st, a_p) &\leq |u| + |v| \\ &\leq (q-p) + (j-i+1) \\ &= (q-i+1) + j-p \\ &\leq j-p \end{aligned}$$

となる。

(QED)

以上より、定理1が証明できれば完全性が証明できることがわかった。それでは、定理1の証明に戻ろう。そのためにまず次の補題2を証明する。これは、

図2. 10のアルゴリズムで用いた関数  $\phi$  を特徴付けるものである.

補題2.  $u \in \text{SUF}(K) - \{\varepsilon\}$ ,  $s = \text{state}(u)$  とする. 集合

$$(\text{PRE}(u) - \{u\}) \cap \text{SUF}(K)$$

の要素のうち長さが最大のもを  $v$  とすれば,

$$\phi(s) = \text{state}(v)$$

となる.

証明  $|u|$  に関する帰納法により証明する.

$|u| = 1$  のとき,  $\phi(s) = 0$ ,  $(\text{PRE}(u) - \{u\}) \cap \text{SUF}(K) = \{\varepsilon\}$  であるから明らかである.

$|u| > 1$  のとき,  $u = ax$  ( $a \in \Sigma$ ,  $x \in \text{SUF}(K)$ ) とする. 集合

$$(\text{PRE}(u) - \{u\}) \cap \text{SUF}(K)$$

の要素を長さの長い順に,  $x_1, x_2, \dots, x_n = \varepsilon$  とし,

$$r_0 = \text{state}(x),$$

$$r_i = \text{state}(x_i) \quad (i = 1, 2, \dots, n)$$

とすれば, 帰納法の仮定により,

$$\phi(r_i) = r_{i+1} \quad (i = 0, 1, 2, \dots, n-1)$$

となる.

(1)  $g(r_m, a) \neq \text{fail}$  となる  $m$  ( $1 \leq m < n$ ) が存在するとき,  $m$  をそのような最小のインデックスとすれば,

$$\phi(r_i) = g(r_m, a) = \text{state}(ax_m),$$

$$ax_m \in (\text{PRE}(u) - \{u\}) \cap \text{SUF}(K)$$

となる.  $m$  のとり方から,  $ax_m$  は  $(\text{PRE}(u) - \{u\}) \cap \text{SUF}(K)$  の要素のうち長さが最大のものであることがわかる.

(2) 任意の  $i$  ( $1 \leq i < n$ ) に対して,  $g(r_m, a) = \text{fail}$  であるとき,

$$\phi(s) = 0, \quad (\text{PRE}(u) - \{u\}) \cap \text{SUF}(K) = \{\varepsilon\}$$

であるから明らかである.

(QED)

補題2から直ちに次の系が得られる.

系.  $u, v \in \text{SUF}(K)$ ,  $u \neq \varepsilon$  であるとき次が成立する.

$$v \in \text{PRE}(u) - \{u\} \Leftrightarrow \text{state}(u) \vdash_{\phi}^+ \text{state}(v)$$

補題3.  $au \notin \text{SUF}(K)$  なる任意の対  $(u, a) \in \text{SUF}(K) \times \Sigma$  に対して,

$$H(u, a) = \{v \in \Sigma^+ \mid auv \in \text{SUF}(K)\}$$

とおく. このとき,  $st = \text{state}(u)$  ならば,

$$\text{shift1}(st, a) = \begin{cases} \min\{|u+v| \mid v \in H(u, a)\} & (H(u, a) \neq \emptyset \text{ のとき}) \\ \text{large (undefined)} & (\text{そうでないとき}) \end{cases}$$

が成立する.

(関数  $\text{shift1}$  は図2. 10のアルゴリズムの関数  $\text{adm}$  で更新される前の (つまり, 図2. 10のAで更新される) 関数  $f$  を表し, 図2. 8(a)の場合のシフト量に相当している. 図3. 1参照)

証明  $g(st, a) = \text{fail}$  (または0) となる各  $(st, a)$  に対し,

$$A(st, a) = \left\{ r \in S \left\{ \begin{array}{l} r \neq 0, g(r, a) \neq \text{fail} \\ \exists r_0, \dots, \exists r_n \in S : \\ r_0 = r, r_n = st, \\ r_{i+1} = \phi(r_i) \quad (0 \leq \forall i < n), \\ g(r_i, a) = \text{fail} \quad (1 \leq \forall i < n) \end{array} \right. \right\}$$

とおく. 図2. 10のアルゴリズムから

$$\text{shift1}(st, a) = \begin{cases} \min\{\text{depth}(r) \mid r \in A(st, a)\} & (A(st, a) \neq \emptyset \text{ のとき}) \\ \text{large} & (\text{そうでないとき}) \end{cases}$$

であることは明らかである. また,

$$A(st, a) \subseteq \{\text{state}(uv) \mid v \in H(u, a)\} \text{ なので,}$$

$$H(u, a) = \emptyset \Rightarrow A(st, a) = \emptyset$$

は明らかである. したがって,

$$\min\{|u+v| \mid v \in H(u, a)\} = \min\{\text{depth}(r) \mid r \in A(st, a)\} \quad (H(u, a) \neq \emptyset \text{ のとき})$$

が成り立つことを示せばよいことになる.

$$|v| = \min\{|z| \mid z \in H(u, a)\}$$

となる  $v \in H(u, a)$  をとり,  $r = \text{state}(uv)$  とおく. このとき,

$$\exists r_0, \dots, \exists r_n \in S : [r_0 = r, r_n = st, r_{i+1} = \phi(r_i) \quad (0 \leq \forall i < n)]$$

である.  $r \in A(st, a)$  であることを示したい. そのためには,

$$1 \leq \forall i < n : g(r_i, a) = \text{fail}$$



を示せばよい。そこで、

$$1 \leq \exists j < n : g(r_j, a) \neq \text{fail}$$

と仮定する。このとき、 $\text{state}(uv') = r_j$  なる  $v' \in \text{SUF}(K) - \{\varepsilon\}$  をとると

$$v' \in H(u, a), |v'| < |v|$$

となり、矛盾する。これによって、

$$A(st, a) \neq \emptyset,$$

$$\min\{|u+v| \mid v \in H(u, a)\} \geq \min\{\text{depth}(r) \mid r \in A(st, a)\}$$

が示された。逆向きの不等号は明らかである。

(QED)

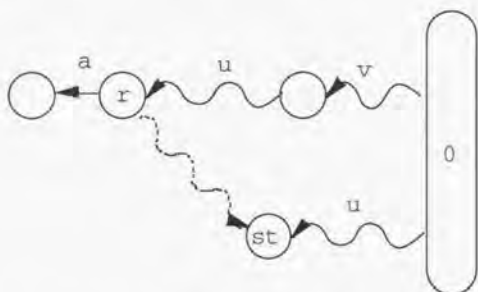


図 3. 1 *shift1* 関数の構成

次の補題は、図 2. 10 のアルゴリズムの関数 *ddm* を特徴付けるものである。

補題 4.  $\text{SUF}(K)$  の空語  $\varepsilon$  でない任意の要素  $u$  に対して、

$$V(s) = \{v \in \Sigma^+ \mid \exists y \in \text{SUF}(u) : yv \in K\}$$

とおく。ただし、 $s = \text{state}(u)$  である。このとき、

$$\text{ddm}(s) = \min\{|uv| \mid v \in V(s)\}$$

が成り立つ。

(この関数は図 2. 8 (b) の場合のシフト量に相当している。 図 3. 2 参照)

証明 任意の状態  $s$  について、

$$B(s) = \{r \in \text{state}(K) \mid r \vdash_{\neq}^+ s\}$$

とおくと、

$$\text{ddm}(s) = \begin{cases} \text{minlen} & \text{if } s = 0 \\ \min\{\{\text{depth}(r) \mid r \in B(s)\} \cup \{\text{ddm}(\text{fath}(s)) + 1\}\} & \text{if } s > 0 \end{cases}$$

となる。ただし、 $\text{fath}(s)$  は  $s$  の親の状態を示す。

$s=0$  のとき、成り立つことは明らかである。  $s>0$  のとき、  $s' = \text{fath}(s)$  とすれば、

$$V(s) = \{v \in \Sigma^+ \mid uv \in K\} \cup V(s')$$

となるので、この場合にも成り立つ。ゆえに、補題は明らかである。

(QED)

補題3と補題4が成り立てば、定理1が成り立つことは自明である。これで、完全性も証明されたことになる。

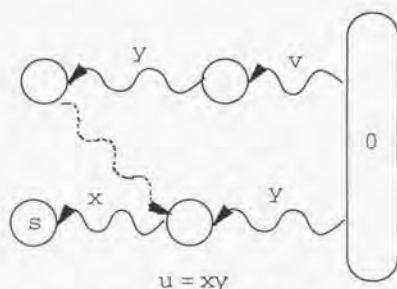


図3. 2 ddm 関数の構成

### 3. 4 検索実験

3. 2で述べたようにFAST法の効率は文字種の数  $q$ 、パターン長  $m$ 、パターン数  $k$  の組み合わせによって大きく変化する。そこで  $q$ 、  $m$ 、  $k$  をいろいろに変えて、FAST法による文字列照合の実験を行った。実験には一様乱数で発生された1Mバイトのテキストを3種、パターンを16種使用した。  $q$  と  $m$  の組み合わせに対して、  $k$  を変化させて効率 (probe rate) を測定した結果を図3. 3に実線で示す。図中には推定式 (3. 2の(2)式) から得られる計算値も一点鎖線で示す。近似式 (3. 2の(3)式) による計算値も破線で示す。AC法では効率は常に1であるので、図中の二点鎖線より下ではAC法より効率が良いことになる。  $q=2$  の場合を除いて推定式と実験結果がよく一致することがわかる。  $q=2$  のとき誤差が大きくなるのは、3. 2で述べたように推定式を求める際に単純化のために近似を入れているため、

この影響は  $q$  が小さいほど大きくなる。近似式の方も  $q$  が大きく  $k$  が小さい時にはよく一致することがわかる。

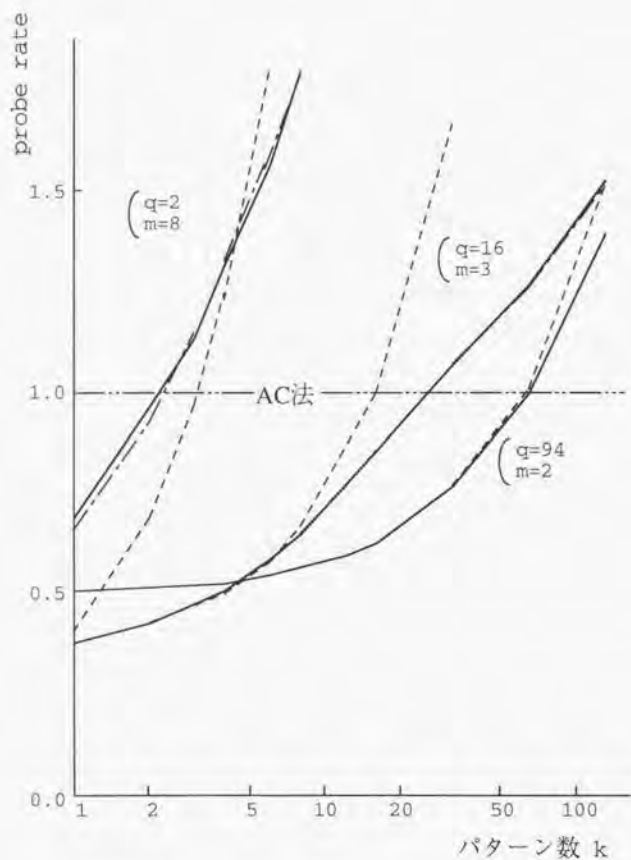


図 3. 3 FAST法の probe rate

図3. 4には、推定式から求められる probe rate が1となる、 $k$  と  $m$  との関係をも  $q$  をパラメータとして示す、実線より下の領域でAC法より効率が良いことになる。 $q$  がかなり小さい場合（例えば  $q=5$ ）でもパターン長が長くなるにつれAC法より有利となる領域が急速に拡大することがわかる。

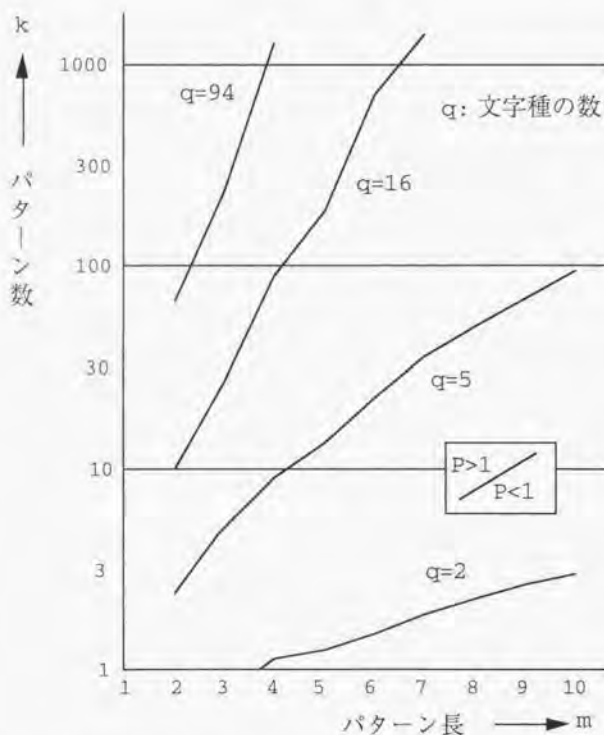


図3. 4 probe rateが1となるパターン長( $m$ )とパターン数( $k$ )の関係

ところで、probe rate は文字列が探索される割合を示したものに過ぎないので、実際に計算機上で実行した場合に上記の効率が得られるかは疑問である。2. 4 で述べたように Horspool は BM 法に probe rate が下がるような改修を行って、実行時間の向上を図っている。そこで、C 言語でコーディングして、AC 法と FAST 法の実際の実行時間の比較を行った。図 3. 5、図 3. 6 にその結果を示す。AC 法で 1 個のパターンを検索したときを 1 とし、それとの比で実行時間を示している。ただし、pmm を作成するための前処理時間やディスクからデータを読み取るための時間は除いている。前者はデータ量が大きくなれば、無視できるものであるし、データをメモリー上に転送する時間はハードウェア構成に異存するだけで AC 法と FAST 法では変わらないからである。つまり、データがメモリー上にあるという条件で探索時間のみを測定した。使用した計算機は Sun Sparc Station 5 である。

図 3. 5、図 3. 6 を見ると probe rate と実測値が良く一致していることが分かる。特に、パターンが少ない場合には誤差が極めて少ない。パターン数が増加するにつれて probe rate より大きな値を取るようになるが、これはパターンを発見したときに出現パターンの種類と出現位置などを記録する処理時間が必要となるからである。AC 法の方も同様で、2 点鎖線で示すように実際にはパターン数が増加すると探索時間がパターンの記憶処理の分だけ増加している。各パターン数毎に AC 法との比を求めたものを破線で示してあるが、この線は probe rate よりも下にくることが分かる。つまり、AC 法との比較でいうなら図 3. 4 で示したよりも若干 FAST 法の方が有利な領域が広いことが想像される。

実際の探索時間であるが、Sun Sparc Station 5 で 1 MB のランダムデータを探索するのに、文字種 16、パターン長 3 の場合、パターン 1 個で 0.158 秒、8 個で 0.264 秒、64 個で 0.580 秒であった。実際の探索例で考えて見ると、AP 電 1 年分の記事（約 250 MB）から、8 個のキーワードを探索するなら（データをディスクから転送する時間を除いて）約 66 秒で済むことが分かる。（英語の文字の単純頻度を考慮したエントロピーは 4 ビット程度と言われている。8 個のキーワードの最短の長さはせいぜい 3 くらいだと思われるから、文字種 16、パターン長 3、パターン数 8 と考えても大差はない。もし、最短のパターン長が 4 なら約 53 秒になる。）これは、通常の用途には十分に実用的な速度であるということが出来る。

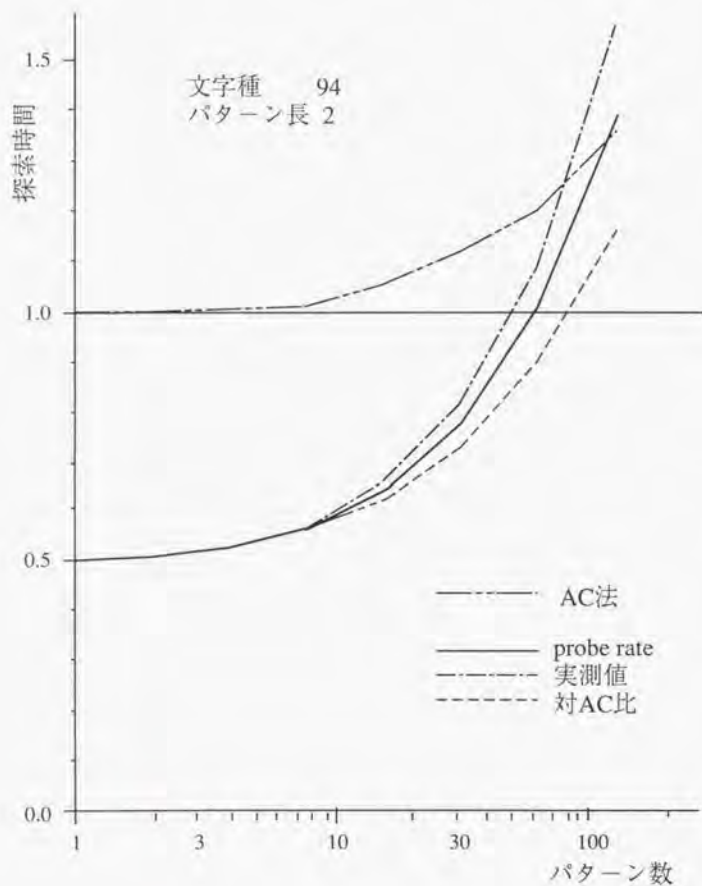


図 3. 5 実際の探索時間 (文字種 94、パターン長 2 の場合)

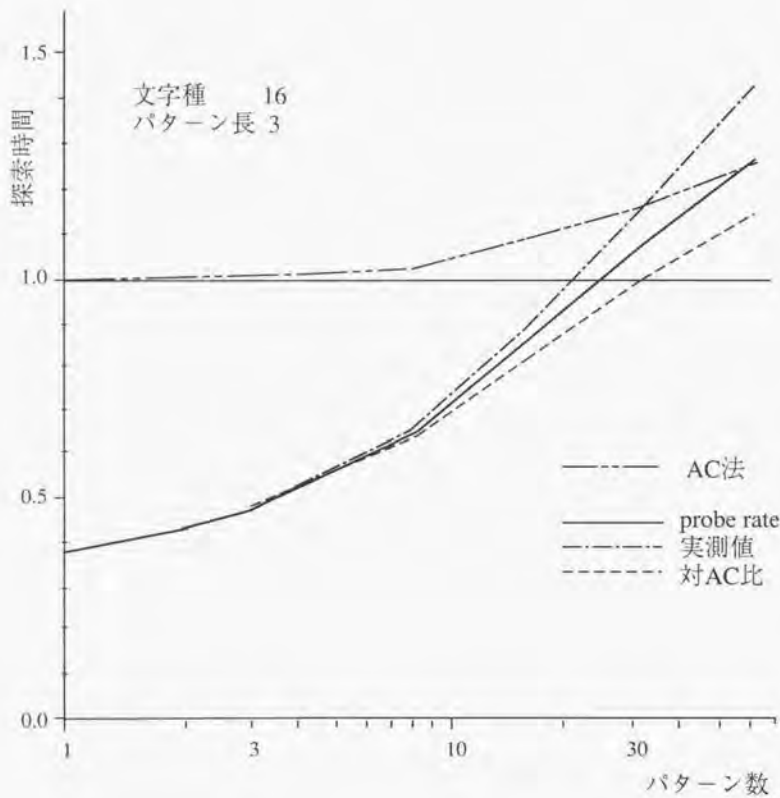


図 3. 6 実際の探索時間 (文字種 16、パターン長 3 の場合)

さて、直感的には、FAST法の probe rate は、 $m$  や  $k$  に対して単調増加、 $q$  に対しては単調減少するように思われる。しかし、 $q$  が大きく、パターン長が長いときには  $k$  に対して単調増加にはならない。この1例 ( $q=94, m=10$  のとき) を図3.7に示す。実測値(実験規模は図3.3と同様)は×印で示すが、推定値とよく一致していることがわかる。これを見ると、パターン数に関して探索効率が単調増加でないこと、つまりパターン数が増えると反対に効率が増す部分が存在することがわかる。この理由は定性的には次のように考えられる。

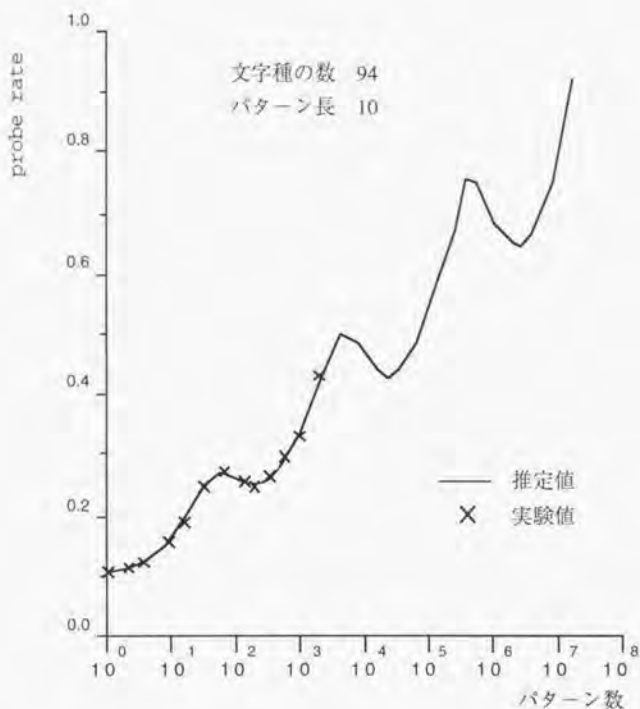


図3.7 長パターン時のFAST法のprobe rate



一般にパターン数が増すと最初の照合に成功する可能性が多くなる（あるいは失敗してもパターン群のどこかと一致するのでシフト量が小さくなる）ので効率が低下する。パターンがある程度以上多くなると最初の照合で失敗する可能性は極めて小さくなる。ところが、この場合2文字連続して照合に成功する可能性は小さいので、照合に失敗したときのパターン群のシフト量は（パターン群のどこかと2文字連続して一致する可能性は小さいので）大きくなる。この結果、最初の照合で失敗するがパターン群のどこかと一致するために少ししかパターン群をシフトできない場合に比べて、最初の照合に成功するけれども次の照合（2文字分の照合）に失敗するために大きなシフトが起きるときの方が平均的に照合回数が少なくなることが起こり得るのである。もちろん、パターン数がさらに増加すると2文字連続の照合成功可能性が増してくるので再び効率は低下する。より詳細に述べると、推定式

(3. 2の(2)式)の分子(照合回数の期待値)は単調増加となるが、分母(シフト量の期待値)はそうならない。例えば、図3. 7 ( $q=94, m=10$ )の場合、パターン数が57, 186, 4638のとき、分子はそれぞれ1.4649, 1.8849, 2.4161となり分母はそれぞれ5.4297, 7.4253, 4.8735となる。

この現象はパターンが長いほど、また文字種が多いほど起こりやすくなる。例えば  $q=9$  の時には  $m \geq 17$  でこの現象が起こるが、 $q=57$  なら  $m \geq 6$  で起こる。また、パターンが非常に長い場合には、図3. 7のようにこの現象が2回以上起こることがある。これは、 $j$ 回の照合に成功しやすく、 $j+1$ 回の照合は成功しにくいとき（つまり  $t_j \gg t_{j+1}$  のとき）に起こりうるのである。

この現象は、例えば図3. 7でprobe rateの曲線が山となるようなパターン数の場合、ダミーのパターンを加え谷の部分に移動させることにより、効率の改善が図れることを示唆している。実際にそうなりうることを例で示す。図3. 7で、最初の山はパターン数が57のところ、このときのprobe rateは0.2698である。最初の谷はパターン数が186のところ、probe rateは0.2539である。したがって、もし探したいパターンの数が57付近であれば、ダミーのパターンを130程度加えて探索すれば0.016近くprobe rateが低下することになる。もちろん、この場合パターンが検出された後で、本当のパターンかダミーパターンかを判定する処理が必要となる。このためには本当のパターンだけで作ったgoto関数(パターン照合機械)を用意しておいて、検出されたパターンをもう1度これで照合すればよい。これは、等価的なprobe rateの増加を招くが、この例の場合、パターン数が186のときの、パターン検

出確率はわずか $3.47 \times 10^{-18}$ であり、最悪でも1回当たりパターン長分しか検出の手間  
 はかからないから、全く問題にはならない。ただし、以上の考察は3章での前提  
 (つまりテキストもパターンもランダム)に基づいていることに注意を要する。

### 3.5 メモリー削減手段<sup>4)</sup>

F A S T法の欠点はgoto関数の記憶に多大のメモリー容量と前処理時間を必要と  
 する点である。例えば、goto関数1個当たり2バイトのメモリーを使用すれば、  
 p m mの状態数がsで文字種がqの場合 $2 s q$ バイトのメモリーが必要である。パ  
 ターン数をk、平均のパターン長をmとすれば、sは大体 $k m$ 程度になる。したがっ  
 て文字種が少なくパターンも少ないときには記憶容量もあまり大きくならない

( $q=53, k=10, m=3$ なら約3 Kバイト)が、文字種が多くパターン数も多いとき  
 には莫大な記憶容量が必要となる。例えば、J I S漢字コードを扱うときにはqは約  
 7000になり、 $k=1000, m=3$ なら約4.2 Mバイトも必要である。この場合単に表が大  
 きくなるばかりでなく、表の初期化等の前処理時間も長くなってしまふ。

有川らはこの様なときにこれを改善する方法を提案している。一つは“ $2 \times n$ 実  
 現”法<sup>5)</sup>を用いる方法である。パターンが第2章と同じ“state”, “east”,  
 “smart”の場合の例を表3.1に示す。これはp m mを $2 \times n$ の表で実現する方

表3.1 p m mの $2 \times n$ 実現

#### (a) A C法の場合

state	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
input	!	s	t	a	t	e	?	m	a	r	t	?	e	a	s	t	?
failure	/	12	7	1	1	1	13	1	1	1	1	1	0	1	1	7	3

#### (b) F A S T法の場合

state	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
input	!	e	t	a	t	s	?	t	s	a	e	?	r	a	m	s	?
failure	/	7	1	8	1	8	9	0	12	1	1	2	1	1	1	1	1

法で、goto関数は  $\text{goto}(s, \text{input}) = s+1$  として復元できる。この表においては、初期状態は1であり、!は任意の文字を、/は無定義を、?は計算機で使われない文字を示している。この方法を用いれば、画期的にメモリーを削減することが可能であるが、この方法では表の参照回数が多くなってFAST法の利点が著しく損なわれてしまう。参照回数が多い状態0(表3. 1では状態1に相当)からの遷移だけを表2. 1の形で残し、それ以外を“ $2 \times n$ 実現”法で行えば、この欠点はかなり救済されると思われる。

有川らは別のメモリー節減手段も提案している<sup>6)</sup>。この方法は文字コードを2つの部分コードに分割し、1回の遷移を2回の連続した遷移に置き換える方法で、“文字符号分割による実現”と有川らは呼んでいる。パターン照合機械(pmm)をこの方法で表現した例をもとのpmmとともに図3. 8に示す。もとのパターンは“AC”, “BAA”, “BB”であり、文字コードはそれぞれ部分コードで

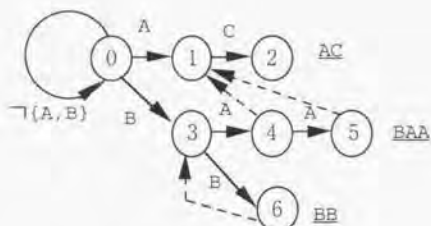
A=aa, B=ab, C=ba

で表されると仮定している。図3. 8で大きな円は主状態を小さな円は中間状態を示している。主状態と中間状態を分けて扱っているのは誤検出を防ぐためであり、主状態からのfailure遷移の相手は主状態でなければならず、中間状態からのfailure遷移の相手は中間状態でなければならない。破線はfailure遷移を示しているが、主状態2, 4, 5から状態0への遷移、8以外の中間状態から状態11への遷移は省略してある。

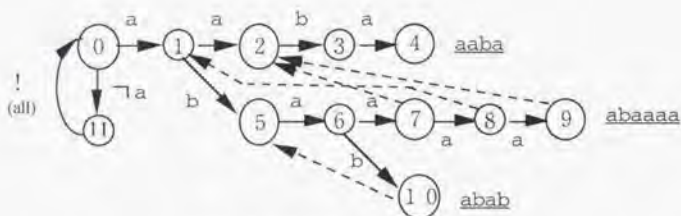
この方法はFAST法にとってより現実的だと考えられる。例えば、JIS漢字コードの2バイトを上位バイトと下位バイトに分割すれば、文字種は94で済むことになる。前と同様  $k=1000$ ,  $m=3(\times 2)$  を仮定すれば、約1.1Mバイトですむことになる。つまりこの場合の記憶容量は元の約  $2/\sqrt{q}$  倍に減少する。ただし、その様な分割符号に対してFAST法を直接適用すれば1バイトずれた誤った検出が起る可能性がある。これを回避するためには、テキストおよびパターン中の文字を牛島らの方法<sup>7)</sup>等によってすべて同じ符号長に変換しておいてから、上述したように中間状態を導入して、上位バイトから下位バイトへの遷移と下位バイトから上位バイトへの遷移を分離すればよい。あるいは牛島らが行なったように、パターン検出後にテキスト上の位置情報からずれ読みか否かを判定すればよい。

メモリーと効率の間にはトレードオフが存在する。例えば、文字種が少なければ文字をまとめて扱うことによって、走査回数を減らすことによって効率化できる。

この場合、当然のことだが、文字種が増えるのでパターンを保持するためのメモリー容量が増大する。その逆は上述した通りである。文字種が多いときは文字コードを分割してメモリーを減らすことができるが、これは等価的にテキスト増大（ $n$ 分割なら $n$ 倍）を招くので、効率が低下する。しかし、FAST法の場合（AC法とは違って）、パターン長は長くなるのでその分効率が良くなり、分割による効率低下の影響は少ない方法であるといえる。



(a) もとの pmm



(b) 文字符号分割による pmm

図 3. 8 パターン照合機械の文字符号分割による実現

### 3. 6 あとがき

本章では、FAST法に関する諸性質、つまり効率、正当性の証明、実験結果、メモリー削減方法について述べた。まず、効率を求めるために、照合の過程をマルコフモデルで近似した。これをもとに文字の照合に成功する確率と文字列のシフト量の期待値を計算し、多元一次方程式を解くことにより、効率 (probe rate) が推定できることを示した。さらに、文字種がパターン長とパターン数の積より十分大きい場合には簡単な式で表せることを示し、 $1/\text{パターン長}$ に近い probe rate で検索が可能であることを証明した。

アルゴリズムが本当に正しく動作するか否かは重大な関心事である。いくら、高速な方法でも、誤検索 (不要なものも出力する、あるいは検索漏れがある) があるのでは何にもならない。FAST法の場合、停止性と健全性はほとんど自明であった。完全性 (検索漏れがないこと) を証明するために、3つの補題と1つの定理を証明し、そこから導きだされる補題を証明することでアルゴリズムの正当性を証明した。

FAST法の効率の良さは直感的には明らかであるが、実際にはどれくらいなのかは興味のあるところである。また、3. 2で示された効率の計算法が実際とどれくらい一致するのも気にかかるところである。実験を行なって、実際に効率の良い検索が可能なこと、効率の推定式は文字種が少ない場合にはよく実験値と一致することを示した。さらに、文字種をパラメータとしてパターン長を横軸として、probe rate が1となるパターン数をグラフで示し、AC法より優れている領域を明らかにした。文字種が少ない場合 (例えば5) でもパターン長が長くなると急速にAC法より優れた領域が拡がるのが明らかになった。このグラフを用いれば、FAST法を使うべきか、AC法を使うべきかをユーザはあらかじめ検討することができる。probe rate と実際にプログラミングしたときの探索時間の差異も気になるところである。実験によって、AC法に対してほとんどprobe rate と一致した効率で検索が可能であることを示した。また、FAST法はパターンが長い場合、直感に反してパターン数に対して効率が単調減少にならないことを示した。さらに、その理由についても考察を加えた。

FAST法の最大の欠点は、goto関数の記憶に多大なメモリーと前処理時間を必要とすることである。言い換えればメモリーを多く消費することで、効率を上げて

いる方法とすることができる。文字種が少ない場合には問題とはならないが、日本語コードなど文字種が多い場合には問題となってくる。メモリー削減手段として2つの方法を上げ、そのうち“文字符号分割による実現”がより現実的な方法であることを示した。また、2バイトコードを扱うときの問題点と対処方法についてもふれた。

F A S T法のもう1つの欠点はワイルドカードが使えないことである。文字列検索では“abc\*xyz”, “table?.1” (\*は任意の文字列; ?は任意の1文字)のように“任意”を指定したいことがよくある。A C法ではどの文字とも一致する状態を作ったり、それを考慮してfailure関数を作成すれば簡単にワイルドカードに対応できるがF A S T法ではそれができない。もともと、完全一致を前提にしてそれを掘り所として効率を上げているアルゴリズムであるからワイルドカードへの対処が困難であることは当然である。しかし、“abc\*や\*xyz”のようにワイルドカードが後ろあるいは前にしかないのなら、“abc”や“xyz”とし、ワイルドカードを無視すれば良いし、“abc?vwxyz”のようにパターン長が長くてワイルドカードが前方にあり(後続文字が短くないという意味)、任意長でなければ、?の部分全文字種に置き換えても(その分パターン数が増えるが)、この章で示した理論や実験結果からそれほどF A S T法の利点は損なわれないことが予想される。

F A S T法は実際にすでに機械翻訳の訳語候補を目標言語のテキストと照合して決定する手法<sup>8)</sup>の中でも使われている。現在、著者は自然言語処理の研究に従事しており、文字列検索が必要な様々な処理にF A S T法を利用していくつもりである。

## 第4章 FAST法の2次元パターン検索への拡張

### 4.1 まえがき

最近、土地利用区分図などのメッシュデータのコンピュータ利用や、デスクトップパブリッシングの発展に伴ない、文字列のような1次元パターンだけでなく、2次元パターンの高速な探索法も求められている。

1次元パターンの探索アルゴリズムを2次元に拡張した例としては、BirdやZhu-TadaokaやBaeza-Yatesの仕事がある。Birdは行方向をAC法で、列方向をKMP法で検索することによって2次元パターンの探索を実現している<sup>1)</sup>。Zhu-Tadaokaは2つの方法を提案している。どちらもテキストをあらかじめパターンの列方向の長さ分の文字列を用いてハッシュ値を計算しておいて、パターンのハッシュ値を1つはKMP法でもう1つはBM法の変形版を用いて検索する方法である<sup>2)</sup>。Baeza-Yatesは $m$ （パターンの列方向の長さ）行毎にCommentz-Walterの方法で行パターンを検索して、それが見つかった場合にはパターンの検出されるエリアを素朴な方法で検索する方法を提案している<sup>3)</sup>。これらの方法はすべて1つの2次元パターンだけを探索対象としていて、パターンが複数の場合は考慮されていない。Zhu-Tadaokaは列方向もAC法を使うことによって、Birdの方法を複数の2次元パターンの探索に拡張する方法を考案している<sup>4)</sup>。この方法はAC法を基にしているため、探索されるデータ空間の大きさに比例した探索時間がかかる。FAST法は前述したように複数の1次元パターンを効率良く探索する方法である。これを2次元パターンの探索に拡張できれば、もっと効率の良い探索が期待できる。この拡張方法を考案し、実験によってその有効性を確認したので、それについて述べる。

### 4.2 FAST法の2次元パターン検索への拡張<sup>5)</sup>

FAST法を2次元に拡張するための基本的な考え方は以下のとおりである。2次元パターンとしては、説明の便宜上、ここでは2次元文字列を例にとる。

今、図4.1のようなパターンを（2次元の）テキストから探索する問題を考える。PT<sub>1</sub>、PT<sub>2</sub>、PT<sub>3</sub>が3つの2次元のパターンである。このパターン探索に2次元FAST法では2段階の探索を行う。まずパターンを列方向と行方向に分解し

て考える。すなわちパターンの各行を1次元パターンと見なしてパターン（行パターン）それぞれに固有の記号（行パターン記号：RP記号）を割り当てる。例えば図4の行パターン“a b c d e”には $X_1$ ，“f g h i j”には $X_2$ ……とRP記号を付ける。このとき、同じ行パターンには同じRP記号が付くようにする。このRP記号を列方向に並べたものを列パターンと呼ぶことにする。つまり、 $PT_1$ の列パターンは $(X_1 X_2 X_3 X_4 X_5)^T$ 、 $PT_2$ の列パターンは $(X_6 X_7 X_8 X_9 X_{10})^T$ 、 $PT_3$ の列パターンは $(X_{11} X_{12} X_{13} X_{14} X_{15})^T$ となる。探索の第一段階では行パターンを（通常の）FAST法で探索する。行パターンが見つければ、列方向にFAST法を用いて列パターンを探索する。列パターンが見つかったときに元の2次元パターンが存在していることがわかる。FAST法と同様にテキストの照合位置を決めるためにgoto関数を用いるが、この作成法はFAST法と同じなので説明を省略する。た

$PT_1$ : abcde = $X_1$   
 fghij = $X_2$   
 hikde = $X_3$   
 demlk = $X_4$   
 bcamj = $X_5$

$PT_2$ : abcd = $X_6$   
 mcag = $X_7$   
 dlfk = $X_8$   
 ghij = $X_9$   
 mcag = $X_7$   
 dlfk = $X_8$

$PT_3$ : kgcde = $X_{10}$   
 demlk = $X_4$   
 bcamj = $X_5$

図4. 1 2次元パターンの一例



だしgoto関数は行パターン，列パターンそれぞれに用意しておく。図4. 2に行パターン用のパターン照合機械の一部を示す。表4. 1に行パターンの，表4. 3に列パターンのgoto関数を示す。表4. 2には行パターンの出力関数を示す。

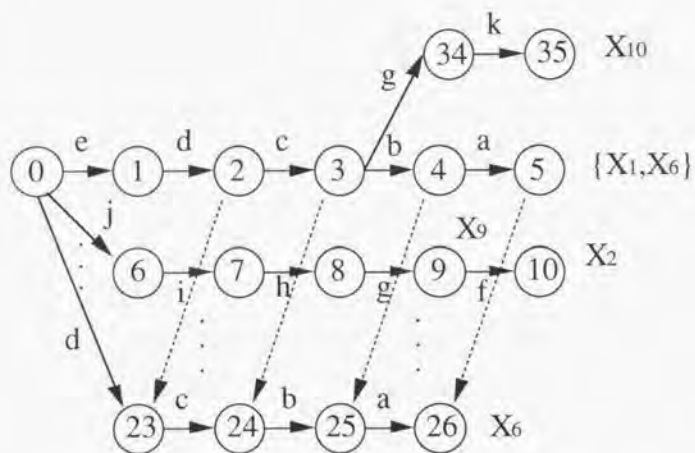


図4. 2 2次元FAST法における  
行パターン用照合機械

表4. 1 行パターンのgoto関数の例

文字種/状態	0	1	2	3	23	24	25	31	32	33	34	35
a	-1	-5	-5	-6	-4	-5	26	-6	-7	-8	-7	-8
b	-2	-5	-5	4	-4	25	-6	-6	-7	-8	-7	-8
c	-1	-5	3	-6	24	-5	-6	-6	-7	-8	-7	-8
d	23	2	-5	-6	-4	-5	-6	-6	33	-8	-7	-8
e	1	-5	-5	-6	-4	-5	-6	-6	-7	-8	-7	-8
f	-1	-5	-5	-6	-4	-5	-6	-6	-7	-8	-7	-8
g	27	-5	-5	34	-4	-3	-6	-6	-7	-8	-7	-8
h	-2	-5	-5	-6	-4	-5	-6	-6	-7	-8	-7	-8
i	-1	-5	-5	-6	-4	-5	-6	-6	-7	-8	-7	-8
j	6	-5	-5	-6	-4	-5	-6	-6	-7	-8	-7	-8
k	14	-5	11	-6	-2	-5	-6	-6	-7	-8	35	-8
l	-1	-5	-5	-6	-4	-5	-6	32	-7	-8	-7	-8
m	-1	-5	-5	-6	-4	-5	-6	-6	-7	-8	-7	-8
(other)	-4	-5	-5	-6	-4	-5	-6	-6	-7	-8	-7	-8

表 4. 2 行パターン出力関数の例

state	rout
5	{X <sub>1</sub> , X <sub>6</sub> }
9	X <sub>9</sub>
10	X <sub>2</sub>
13	X <sub>3</sub>
18	X <sub>4</sub>
22	X <sub>5</sub>
26	X <sub>6</sub>
30	X <sub>7</sub>
33	X <sub>8</sub>
35	X <sub>10</sub>
(other)	0

表 4. 3 列パターンのgoto関数の例

文字種/状態	0	1	2	3	4	5	6	7	8	9	10	11	12
X <sub>1</sub>	-3	-4	-5	-6	5	-8	-4	-5	-6	-7	-8	-9	-6
X <sub>2</sub>	-3	-4	-5	4	-7	-8	-4	-5	-6	-7	-8	-9	-6
X <sub>3</sub>	-2	-4	3	-6	-7	-8	-4	-5	-6	-7	-8	-9	-6
X <sub>4</sub>	-1	2	-5	-6	-7	-8	-4	-5	-6	-7	-8	-9	-6
X <sub>5</sub>	1	-4	-5	-6	-7	-8	-4	-5	-6	-7	-8	-9	-6
X <sub>6</sub>	-3	-4	-5	-6	-7	-8	-4	-5	-6	-7	11	-9	-6
X <sub>7</sub>	-1	-4	-5	-6	-7	-8	7	-5	-6	10	-8	-9	-6
X <sub>8</sub>	6	-4	-5	-6	-7	-8	-4	-5	9	-7	-8	-9	-6
X <sub>9</sub>	-2	-4	-5	-6	-7	-8	-4	8	-6	-7	-8	-9	-6
X <sub>10</sub>	-2	-4	12	-6	-7	-8	-4	-5	-6	-7	-8	-9	-6
(other)	-3	-4	-5	-6	-7	-8	-4	-5	-6	-7	-8	-9	-6

#### 4. 2. 1 2次元FAST法のアルゴリズム

上述の考え方を厳密にアルゴリズム（照合部分のみ）で記述すると、図4. 3のようになる。ただし変数の宣言や関数の引数や初期化の一部は省略してある。  $a_{ij}$  はテキストの  $i$  行  $j$  列目の文字を示している。  $map$  は行パターンの探索結果を保持しておくための2次元配列で、その容量は最大の列パターンを検出できるように、

テキストの行長 ( $m$ )  $\times$  (列パターンの最大長 ( $s$ ) + 1)  
となっている。

-----  
テキストの文字列  $x$  を  $x = a_{11}a_{12}\dots\dots a_{1n}$  とする。

$\dots\dots a_{ij}\dots\dots$   
 $a_{m1}a_{m2}\dots\dots a_{mn}$

すなわち  $m \times n$  の2次元データである。

行パターンのgoto関数を  $rg$ 、列パターンのgoto関数を  $cg$ 、

出力関数を  $cout$  とする。行パターンに対する出力関数

(つまり列パターン探索の入力) は  $rout$  である。

出力は  $x$  中出现するパターンとその位置である。

$length$  は行パターンの長さを求める関数である。

$map$  は  $mm$  行 (現在の探索対象行に相当) は0で、その他の行は#で初期化されているものとする。

```
1 begin
2   qc:=cmin; /* cmin は列パターンの最小長 */
3   while qc ≤ n do
4     begin
5       mn:=cmax; /* cmax は列パターンの最大長 */
6       rst:=0;
7       dm:=cmin;
8       qr:=rmin; /* rmin は行パターンの最小長 */
9       while qr ≤ m do
10        begin
```

```

11     if rg (rst, aqc,qr) < 0 then
12         begin
13             qr:=qr-rg (rst, aqc,qr);
14             rst:=0
15         end
16     else
17         begin
18             rst:=rg (rst, aqc,qr);
19             if rout(rst) > 0 then
20                 begin
21                     if map(mm,qr) = 0 then
22                         map(mm,qr):=rout(rst);
23                     else if map(mm,qr) > 0 then
24                         map(mm,qr):= #;
25                         /* # は "未探索" を示す記号 */
26                     end
27                     cst:=cg (0, rout(rst));
28                     if cst < 0 then
29                         dm:= min(dm, -cst);
30                     else
31                         begin /* 列パターンの探索 */
32                             len:=length (rout(rst));
33                             dm:=csearch (cst,mm,dm,len);
34                         end
35                     end
36                     qr:=qr-1;
37                 end
38             end
39             qc:=qc + dm;
40             maprar (map, dm) /* map の並べ直しと初期化 */
41         end
42     end
43 end

```

```

41 function csearch (cst,mm,dm,len)
42 begin
43   mc:=mm; qs:=qc-1;
44   while mc ≤ cmax do
45     begin
46       if cout(cst) ≠ empty then
47         print qr,qc-clen(cst)+1,cout(cst);
48         /* clen(cst) は cst に対応する列パターンの長さ */
49         mc:=mc-1;
50         if map(mc,qr) = # then
51           begin /* map の部分計算 */
52             j:=qr+len-1;
53             tst:=0;
54             while j ≤ (qr+len-1) do
55               begin
56                 tst:=rg (tst, aqs,j);
57                 if tst < 0 then return(dm);
58                 else
59                   begin
60                     if j = qr and rout(tst) > 0 then
61                       begin
62                         if cg (cst,cout(tst)) < 0 then
63                           return(min(dm,mc-cst-cmax));
64                         else break
65                       end
66                     j := j-1;
67                     if j < qr then return(dm)
68                   end
69                 end
70             end
71           end
72         end
73       end

```

```

70     else
71         begin /* map 上の探索 */
72             if eq (cst,map(mc,qr)) < 0 then
73                 return(min(dm,mc-cst-cmax));
74             end
75             qs := qs - 1;
76         end
77     end

```

図4. 3 2次元FAST法のパターン照合アルゴリズム

2次元FAST法ではまずテキストの  $cmin$  行目 ( $cmin$  は列パターンの最小の長さで図4. 1の場合には3) から行パターン探索 (図4. 3の8行~35行) を始める。何故なら、もし  $cmin$  行目に2次元パターンの最下段に当たる行パターン (例では、 $X_3$  と  $X_6$  に相当) が見つからなければ、1~ $cmin$  行の間に2次元パターンが存在するはずがないから、それまでの行を探索対象から外すことによって効率を上げることができるからである。 $cmin$  を列パターンの最小の長さにするのは、最小の列パターンを検出し損なわないようにするためである。したがって、次に探索対象とする行までのスキップ幅 ( $dm$ ) もこれより大きくすることはできない。

行パターンが見つかったときには結果 (RP記号) は  $map$  に書き込まれる (図4. 3の19行~25行)。このとき  $map$  の対応行 ( $mm$ ) は、最大の列パターンも検出するために列パターンの最大長に当たるところにとらなければならない。それが  $qc$  の初期値と  $mm$  の初期値が ( $cmax-cmin$ ) 分ずれている理由である。1行分の検索が終わると、 $qc$  も  $map$  も  $dm$  分シフトされる (図4. 3の37行, 38行)。

$map$  に検索結果を格納する際、すでにRP記号が書き込まれている場合があるので注意が必要である。これは1つの行パターンが他の行パターンの prefix である場合に起こる。例えば図4. 1の "abcde" が検出されるときには必ず

"abcd" も検出される。この場合には  $map$  上に2つのRP記号  $X_1$  と  $X_6$  を記憶しなければならぬが、2次元FAST法ではこうした場合、RP記号を "未探索"

に変えることで対処している。もちろん、リスト構造などを用いて2つ以上のRP記号を同一場所に記憶できるようにすることも可能であるが、map上の情報はスパースなので記憶容量の節約のためこの方法の方が望ましいと思われる。確かに、“未探索”に変え情報を消滅させてしまうのは非効率だが、このために生じる余分な探索の手間はごくわずかで、問題になるとは思われない。

もう1つの解決手段は小さいパターン( $X_0$ )が大きいパターン( $X_1$ )と同じ位置で検出される場合、小さいパターンを含む列パターンの替わりに大きいパターンを含むものも列パターンとして登録する方法である。すなわち、PT<sub>2</sub>の列パターンとして( $X_0X_7X_8X_0X_7X_8$ )<sup>1</sup>だけでなく、( $X_1X_7X_8X_0X_7X_8$ )<sup>1</sup>も考慮するわけである。こうすれば、探索すべき列パターンが多くなるという欠点が生ずるが、いずれにせよ $X_1$ が検出されれば $X_0$ も検出されるわけであるから効率の低下は、単に見かけ上のことに過ぎない。map上に2つのRP記号を記録しなければならない場合を考慮する必要がなくなるので、この方法の方が実用的かもしれない。

1行の行パターンの探索が終わった後で次にどの行を探索するかは、どの行パターンが検出されたかによって決まる。例えば行パターン探索で $X_2$ と $X_9$ のパターンが検出されたならば、表3のgoto関数から次の行パターン探索は2行先から始めれば良いことがわかる(図4.3の26行~28行)。もし表3のgoto関数が正ならば列パターンで最後尾となりうるRP記号が見つかったことになるのでmapを対象に列パターン探索を開始する(図4.3のcsearch)。すなわち列を1つずつ戻りながら( $mc \leftarrow mc-1$ )、map上の値によって列パターンの探索を進めれば良い(図4.3の71行~74行)。しかし、前述のように2次元FAST法では行をスキップすることが起こるので、全ての行が探索済みであるとは限らない。この場合mapには“未探索”であることを示す記号が入っているので、この値を求めるために部分的に行パターンの探索を行う。(図4.3の50行~69行)これは、主プログラムの行パターンの探索と本質的に同じだが、探索の範囲はlen文字分だけでよい。

また、上記の説明はパターンやテキストが1バイトコードで表現できることを前提としているが、1次元FAST法と同様な方法(3.5参照)で2バイトコードを扱うことも可能である。



#### 4. 2. 2 検索実験

テキスト中の文字の平均照合回数で2次元FAST法の効率Pを評価すると、図4. 3のアルゴリズムからわかるように最良時には

$$P = 1 / (c_{min} \times r_{min})$$

となる。ここで $c_{min}$  ( $r_{min}$ ) は列 (行) パターンの最小の長さである。

2次元FAST法の効率はテキストやパターンの性質と大きさによって大きく変化する。テキストもパターンも一様乱数によって生成し、実験した結果を図4. 4と図4. 5に示す。図4. 4はパターンの大きさを $3 \times 3$ に固定し、文字種の数を変数とし、パターン数によるprobe rateの変化を示したものである。図4. 5は文字種数を16に固定し、パターンの大きさをパラメータとし、パターン数によるprobe rateの変化を示したものである。ただし、probe rateを単に文字照合回数で計るのは不適當であると考えたので、

(文字照合回数 + mapの照合回数) / テキスト文字数  
によって計算している。

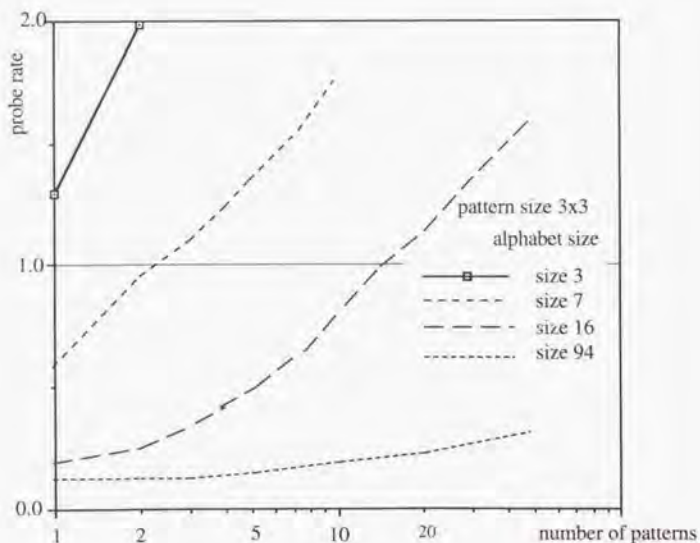


図4. 4 2次元FAST法の探索効率

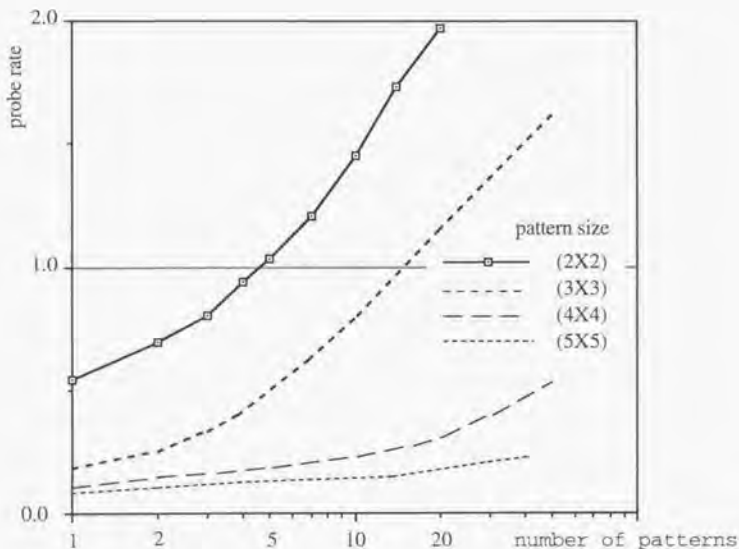


図4.5 文字種16の場合の効率

アルゴリズムから分かるようにテキストの行方向の長さよって列方向のスキップ幅に影響を与えるので効率にも変化する。この実験ではテキストの行方向の幅は512文字で行なっている。

これを見ると当然予想されることながら文字種が多いほど、あるいはパターンが大きいほど2次元FAST法の効率が良いことがわかる。文字種が少なく(例えば16)、パターンが小さい(例えば3×3)場合でもパターン数が少なければ(15以下)、Zhu-Tadaokaの方法(常にprobe rateが1)よりも効率が良いことが分かる。

#### 4. 3 FAST法のn次元パターン検索への拡張

上記のアルゴリズムは簡単にn次元パターンの探索に拡張することができる。パターンを $PT_1[P_{1j}][P_{2j}] \cdots [P_{nj}]$ 、テキストを $TT_1[T_2] \cdots [T_n]$ とする。パターンの第1次元目から第n-1次元目までの配列の要素を固定して考えれば、(重複を含めて) $P_{1j} \times P_{2j} \times \cdots \times P_{(n-1)j}$ 個( $\times$ パターン数)の1次元パターンが得られる。これが、3章で述べた行パターンに相当する。図4.6に3次元パターンの一例を示す。図4.6では $2 \times 2 \times 2$ (このうち3つは重複するので5つ)の1次元パターンが得られる。これにユニークな記号(RP記号)をふり(図4.6の例では $X_1$ から $X_5$ に相当)、次にパターンの第1次元目から第n-2次元目までの配列の要素を固定して考えれば、 $P_{1j} \times P_{2j} \times \cdots \times P_{(n-2)j}$ 個のRP記号を基にした1次元パターンが得られる。(図4.6では $X_1$ から $X_5$ で記述される $Y_1, Y_2, Y_3$ に相当)

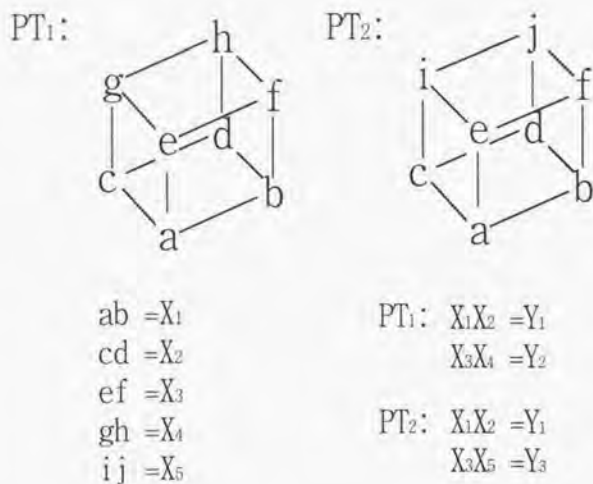


図4.6 3次元パターンの一例

同様にして、1次元ずつ次元が減じたパターンが得られる。これらの1次元パターンに対してFAST法によってgoto関数と出力関数を求める。前節から推測できるように最初の1次元目の出力関数は2次元パターン探索の入力記号となる。順次、 $(n-i-1)$ 次元パターンの出力関数は $(n-i)$ 次元パターン探索の入力記号となる。途中の探索結果を取っておくために、それぞれ、

$$T_1 \times T_2 \times \dots \times T_{(n-1)} \times \max_j (PT_{n_j} + 1),$$

$$T_1 \times T_2 \times \dots \times T_{(n-2)} \times \max_j (PT_{(n-1)_j} + 1),$$

$\dots$ ,  $T_1 \times \max_j (PT_{2_j} + 1)$ , の $(n-1)$ の配列を用意しておけばよい。 $(n-i)$ 次元目のパターンの探索は $T_1 \times T_2 \times \dots \times T_{(n-i-1)} \times \max_j (PT_{(n-i)_j} + 1)$ の $(n-i)$ 次元配列上で行なえばよいことが分かる。

#### 4. 4 あとがき

2次元パターン探索を1次元探索の2つ（行探索と列探索）に分解することにより、FAST法を2次元に拡張するアルゴリズムを考案した。実験によって複数の2次元パターンを同時に高速に探索できることを確認した。この方法はFAST法と同様に文字種が多いほど、パターンが大きいほど効率が良く、パターン数が多くなるにつれて効率が低下することが分かった。

この方法を一般の $n$ 次元パターンの探索に拡張する方法についても述べた。これらの手法は単に文字列の検索にとどまらず、様々な多次元パターンの検索にも利用が可能である。

## 第5章 結論

本論文では、文字列処理やパターン検索の基本である複数文字列の効率的な探索アルゴリズムに関する研究について述べた。情報化社会の進展に伴って、情報検索の重要性がますます増大している。筆者は情報検索のもっとも基礎であるパターン検索を効率化することを目指して検討を行なった。

データベースは必要な情報を、必要な時に利用できるようにするための強力なツールである。データの検索の手掛かりとしては、現在キーワードが最も利用されている。この場合、ユーザに自由にキーワードを指定できるようにするためには全文を対象としてキーワードの出現の有無を調べる手段を提供しなければならない。また、社会には電子化されていてもデータベース化されていない文書データもまだまだ多量に存在している。これを対象に何らかの語句が含まれているか否かを調べるには全文探索よりほかに手段はない。

本論文の主要な目的は、任意の複数の文字列（パターン）を電子化されているデータから効率よく検索する手法を提供することである。

第1章では、研究の背景と目的、特にキーワード検索手法について概観し、インデックス検索手法、長尾・森の方法および全文検索手法の違いと技術要件を明らかにした。合わせて、本論文の意義および位置付けについても明らかにした。

第2章では、全文検索手法のアルゴリズムについて検討した。従来のアルゴリズムのうち、素朴な方法、Knuth-Morris-Prattの方法（KMP法）、Boyer-Mooreの方法（BM法）、Aho-Corasickの方法（AC法）について簡単に説明を行なった。そして、AC法の状態遷移という考え方とBM法の「パターンの右端からのテキストの逆向きの照合」というアイデアを結合することによって、複数パターンを効率良く検索できることを示した。提案した方法（FAST法）の考え方とアルゴリズムについて詳述した。他の複数パターン検索法よりも効率が良いことを明らかにし、パターンが1つだけの場合でもBM法よりも効率が良いことを示した。

第3章では、第2章で提案したFAST法の効率について数学的に検討した。そして、文字種（ $q$ ）がパターンの長さ（ $m$ ）、パターン数（ $k$ ）に対して十分大きい

( $q$ ) $\gg mk$ ) ときには probe rate が  $1/m + (m+1)k/2mq$  で近似できることを示した、ついで、アルゴリズムの正当性を厳密に証明し、このアルゴリズムが信頼に足ることを示した。実験を行ない、効率の推定式がこれとよく一致することを示し、さらに多パターン時には効率が単調に減少しないという現象が起きることについて説明した。FAST法の欠点は他のアルゴリズムよりメモリーを余計に必要とすることである。文字種とパターンの長さとの積が小さい場合は問題にはならないが、JIS漢字コードなどの2バイトコードを対象にする場合には利用が困難になる場合も生じる。こういう場合でも、FAST法を利用できるようにするために“文字符号分割による実現”などのメモリー節減手段を示した。

第4章では、2次元パターン探索の従来のアルゴリズムについて概観した後、パターンを列方向と行方向に分解し、行パターンと列パターンという概念を導入することで、FAST法を2次元パターン探索に拡張する方法について述べた。このアルゴリズムについて詳述し、実験により従来の手法より効率が良いことを示した。さらに同様な発想でFAST法を一般のn次元パターンの検索に拡張する方法についても述べた。

情報化社会の進展はデータ情報の指数関数的に増大し、個人に供給される1日分の情報さえ、24時間を費やしてもはやその全てにアクセスすることができなくなっていると言われている。これに対処するため使い勝手の良い高効率な情報検索手法が強く求められている。文字列の検索は、指定した文字列(パターン)が文書中のどこに現われているかを調べるという単純な作業であるが、電子化された文書の処理にとっては全体の効率に大きく関与する基本操作なので、現在でも重要な研究テーマである。本研究では、有限オートマトンによるパターンの記述と「右端からの照合」というアイデアを結合した新しい全文検索手法について検討し、実験も行ない、その有効性、実現性について検証した。

今後、本研究の概念を生かしながら、様々な情報検索の実際の応用の上で本アルゴリズムの利用を図っていきたく考えている。また、2次元FAST法(あるいはn次元FAST法)については、今後はさらに理論的および実験的に考察を進めて、性質を明らかにするとともに、実用面でも利用を図っていきたく考えている。

## 謝 辞

本論文の作成にあたり、終始御指導、御鞭撻を賜りました東京大学工学部電子情報工学科羽鳥光俊教授に深謝いたします。

東京大学工学部電子工学科齋藤忠男教授、電気工学科田中英彦教授、生産技術研究所今井秀樹教授、電子情報工学科相田仁助教授、相澤清晴助教授には、適切な御指導と御助言を賜りました。深謝します。

本研究はNHK放送技術研究所において行ったものである。本研究を開始するきっかけを与えていただいた江原主任研究員に感謝いたします。本研究を進めるにあたり御指導いただいた沓沢淳之助元画像研究部部长（現、シャープ）、相沢輝昭前研究主幹（現、広島市立大学）、二宮佑一研究主幹に感謝いたします。

有益な御助言をいただいた九州大学理学部有川節夫教授、NHK放送技術研究所の榎並和雅先端制作技術研究部部长、藤原正雄元副部长、井上誠喜氏（現、ATR）、加藤直人研究員（現、ATR）に感謝いたします。

第3章で述べた正当性の証明を研究してくださった九州大学工学部竹田正幸助教授（当時大学院生）に感謝いたします。

3年間のATR出向の間、研究でお世話になったATR自動翻訳電話研究所樽松明前社長（現、電気通信大学）、ATR音声翻訳通信研究所山崎泰弘社長、森元逞第4研究室長、谷戸文廣主幹研究員（現、KDD）、田代敏久研究員（現、マイクロソフト）を初めとする自動翻訳電話研究所に在籍された方々に感謝いたします。また、日頃御討論いただく、畑田のぶ子研究員、金淵培研究員、田中英輝研究員、熊野正研究員、柴田正啓研究員を初めとする画像研究部各位に感謝いたします。

本研究の機会を与えていただいたNHK放送技術研究所杉本昌穂元所長（現、バイオニア）、泉武博前所長、西澤台次所長、中林克己前次長、山田宰次長、大島英男次長に深謝いたします。

## 参考文献

### 第1章

- (1) J. Aoe : "Computer Algorithm -Key Search Strategies-",  
IEEE Computer Society Press (1991)
- (2) 有川節夫, 篠原武 : " 文字列パターン照合アルゴリズム" ,  
コンピュータソフトウェア, Vol.4, No.2, pp.2-23 (1983)
- (3) Boyer,R.S. and Moore,J.S. : "A Fast Searching Algorithm",  
Comm.ACM, Vol.20, pp.762-772 (1977)
- (4) Aho,A.V. and Corasick,M.J. : "Efficient String Matching: An Aid to Bibliographic  
Search", Comm.ACM, Vol.18, No.6, pp.333-340 (1975)
- (5) Kowalski,G.,Meltzer,A. : "New Multi-Term High Speed Text Search Algorithms",  
First International Conference on Computer and Applications (CAT.No.84ch2039-6),  
pp.514-522, XIV+905(1984)
- (6) Commentz-Walter,B. : "A String Matching Algorithm Fast on the Average",  
Lecture Notes in Computer Science, Vol.71, Springer, pp.118-132 (1979)
- (7) 吉村賢治, 日高達, 吉田将 : " 日本語科学技術文における専門用語の自動  
抽出システム" , 情報処理学会論文誌, Vol.27, No.1, pp.33-40 (1986)
- (8) 木本晴夫 : " 日本語新聞記事からのキーワード自動抽出と重要度評価" ,  
電子情報通信学会論文誌D-I, Vol.J74-D-I, No.9, pp.556-566 (1991)
- (9) Nagao M. and Mori S. : "A New Method of N-gram Statistics for Large Number of  
n and Automatic Extraction of Words and Phrases from Large Text Data of Japanese",  
COLING94, pp.611-615(1992)
- (10) Roberts C.S. : "Patial-Match Retrieval via the Method of Superimposed Codes",  
Proceedings of the IEEE, Vol.67, No.12, (1979)
- (11) 有田 健, 津田和彦, 入口浩一, 青江順一 : "複数キーワードによる検  
索の一高速化手法" , 情報処理学会研究会資料, NL104-7 (1994)
- (12) Burkowski F. J. : "A Hardware Hashing Scheme in the Design of a Multiterm  
String Comparator", IEEE Trans. Comput, COM-31, 9, pp.825-834 (1982)



- (13) 井上 潮, 速水治夫, 福岡秀樹, 鈴木健司, 松永俊雄:  
“データベースプロセッサ R I N D A の設計と実現”,  
情報処理学会論文誌, Vol.31, No.3, pp.373-380 (1990)

## 第2章

- (1) Knuth D. E., Morris J. H. and Pratt V. R. : "Fast Pattern Matching in Strings",  
SIAMJ. Comput, 6, 2, pp.322-350 (1977)
- (2) Boyer R.S. and Moore J.S. : "A Fast Searching Algorithm",  
Comm.ACM, Vol.20, pp.762-772 (1977)
- (3) Galil Z : "On Improving the Worst Case Running Time of Boyer-Moore String  
Searching Algorithm", Comm.ACM, Vol.22, No.9, pp.505-508 (1979)
- (4) Aho A.V. and Corasick M.J. : "Efficient String Matching: An Aid to Bibliographic  
Search", Comm.ACM, Vol.18, No.6, pp.333-340 (1975)
- (5) Horspool R.N. : "Practical Fast Searching in Strings", Software-Practice and  
Experience, Vol.10, No.6, pp.501-506 (1980)
- (6) Kowalski G. and Meltzer A. : "New Multi-Term High Speed Text Search Algorithms",  
First International Conference on Computer and Applications (CAT.No.84ch2039-6),  
pp.514-522, XIV+905(1984)
- (7) Commentz-Walter B. : "A String Matching Algorithm Fast on the Average",  
Lecture Notes in Computer Science, Vol.71, Springer, pp.118-132 (1979)
- (8) 竹田正幸 : "複数文字列照合に関する浦谷アルゴリズムの正当性について",  
夏のL A シンポジウム予稿, pp.1-26(1988)
- (9) 浦谷則好 : "高速な複数文字列照合アルゴリズム : F A S T",  
情報処理学会論文誌, Vol.30, No.9, pp.1119-1125 (1989)

## 第3章

- (1) 浦谷則好 : "F A S T法の効率の推定と長パターン時のふるまい",  
情報処理学会論文誌, Vol.32, No.9, pp.1073-1079 (1991)

- (2) Uratani N. and Takeda M. : "A Fast String-Searching Algorithm for Multiple Patterns", *Info. Proc. & Manag.*, Vol.29, No.6, pp.775-791 (1993)
- (3) Takeda M. : "A Proof of the Correctness of Uratani's String Searching Algorithm", *Tech. Report RIFIS-TR-CS-7*, Kyushu University (1988)
- (4) 浦谷則好 : "高速な複数文字列照合アルゴリズム", *NHK技研R & D*, No.9, pp.56-68 (1990)
- (5) Arikawa S. : "One-Way Sequential Search Systems and Their Powers", *Bull. Math. Stat.*, Vol.19, No.3-4, pp.69-85 (1981)
- (6) Arikawa S. and Shinohara T. : "A Run-Time Efficient Realization of Aho-Corasick Pattern Matching", *New Generation Comput.*, Vol.2, No.2, pp.171-186 (1984)
- (7) 尹志熙, 高木利久, 牛島和夫 : "5種類のパターン・マッチング手法をC言語の関数で実現する - 第2回日本語テキストの場合", *日経バイト* 1987年9月号, pp.233-243 (1987)
- (8) 加藤直人 : "目標言語のフルテキスト検索による機械翻訳の訳語選択", *信学技報*, NLC93-32, pp.9-14 (1993)

#### 第4章

- (1) Bird R.S. : "Two Dimensional Pattern Matching", *Inf. Process. Letters*, Vol.6, No.5, pp.168-170 (1977)
- (2) Zhu R.F. and Takaoka T. : "A technique for two-dimensional pattern matching", *Communications of ACM*, Vol.32, No.9, pp.1110-1120 (1989)
- (3) Baeza-Yates R. and Regnier M. : "Fast Algorithms for Two Dimensional and Multiple Pattern Matching", *Lect. Notes Comput. Sci.*, Vol.447, pp.332-347 (1990)
- (4) Zhu R.F. and Takaoka T. : "The Extension of the Aho-Corasick Algorithm to Multiple Rectangular Pattern Arrays of Different Sizes and N-Dimensional Cases", *J. Information Processing*, Vol.11, No.4, pp.271-277 (1988)
- (5) 浦谷則好 : "2次元パターン高速探索アルゴリズム", *第37回情報処理学会全国大会論文集*, pp.1541-1542 (1988)

## 本論文関連の対外発表

### 学会誌論文

- (1) 浦谷：“高速な複数文字列照合アルゴリズム：FAST”  
情報処理学会論文誌, Vol.30, No.9, pp.1119-1125(1989)
- (2) 浦谷：“FAST法の効率の推定と長パターン時のふるまい”  
情報処理学会論文誌, Vol.32, No.9, pp.1073-1079(1991)
- (3) N. Uratani and M. Takeda: "A FAST STRING-SEARCHING ALGORITHM FOR MULTIPLE PATTERNS"  
Information Proc. & Manag., Vol.29, No.6, pp.775-791(1993)
- (4) 浦谷：“複数2次元パターンの高速探索アルゴリズム”  
電子情報通信学会に投稿中

### 研究会・国内発表

- (1) 浦谷：“複数文字列照合アルゴリズム”  
情処全大, pp.57-58(1987)
- (2) 浦谷：“複数文字列照合アルゴリズム”  
信学研報, SS87-27, pp.1-7(1988)
- (3) 浦谷：“2次元パターン高速探索アルゴリズム”  
情処全大, pp.1541-1542(1988)

### その他の対外発表

- (1) 浦谷：“高速な複数文字列照合アルゴリズム”  
NHK技研R & D, No.9, 1990-5月号, pp.56-68(1990)

### 出願特許

- (1) 複数パターンの検索方法 (1987.9)
- (2) 複数2次元パターンの検索方法 (1988.9)

## 本論文関連以外の主な対外発表

### 学会誌論文

- (1) 浦谷, 柴田, 野口, 相沢: “静止画検索システムFORKSの試作”  
情報処理学会論文誌, Vol.28, No.7, pp.758-767(1987)

### 研究会等

- (1) 浦谷: “ランドサットMSSデータによる土地利用区分メッシュデータの細分”, 第7回リモートセンシングシンポジウム, pp.81-84(1981)
- (2) 浦谷: “Prolog処理系の試作と日本語処理への応用”  
NHK技研月報, Vol.27, No.10, pp.424-429(1984)
- (3) 浦谷, 柴田, 江原: “日本語ワードプロセッサにおける放送用辞書の利用”  
NHK技研月報, Vol.30, No.9, pp.301-306(1987)
- (4) 柴田, 浦谷, 井上, 野口: “放送用静止画検索システムFORKS”  
TV学技報, Vol.11, No.31, pp.57-62, OPT'87-18(1987)
- (5) 浪岡, 浦谷, 相沢: “ニュース文における深層格抽出法”  
情処研資, NL77-3, pp.1-8(1990)
- (6) 浦谷, 相沢: “英語ニュースの機械翻訳”  
情処研資, NL78-18, pp.137-143(1990)
- (7) 相沢, 浦谷, 田中: “英語ニュースへの機械翻訳システムの適用”  
信学研報, NLC91-20, pp.31-38(1991)
- (8) 浦谷, 相沢, 田中, 加藤, 畑田, 住吉: “英語ニュースの機械翻訳”  
NHK技研R & D, No.14, 1991-8月号 pp.1-12(1991)
- (9) 相沢, 鎌田, 浦谷: “外電経済ニュースの英日機械翻訳: 新しいアプローチ”  
信学研報, NLC91-45, pp.65-72(1991)
- (10) 浦谷, 鈴木, 永田, 森元, 田窪, 定延, 成田:  
“目的指向型会話文解析システムの機能評価法”  
信学研報, NLC92-10, pp.33-40(1991)

- (1 1) 保坂, 竹沢, 浦谷: "対話データベースを使った無助詞名詞句の分析"  
人工知能学会研究会, SIG-SLUD-9203-1, pp.1-7(1992)
- (1 2) 浦谷, 森元, 谷戸: "音声翻訳システム A S U R A の開発"  
情処研資, CH18-7, pp.49-54(1993)
- (1 3) 谷戸, 竹沢, 嵯峨山, 鷹見, Singer, 浦谷, 森元, 樽松:  
"自動翻訳電話国際共同実験", 信学研報, SP93-23, pp.73-80(1993)
- (1 4) 浦谷: "A T R 会話音声データにおける言語情報表記"  
情処研資, 93-SLP-2, pp.1-4(1993)

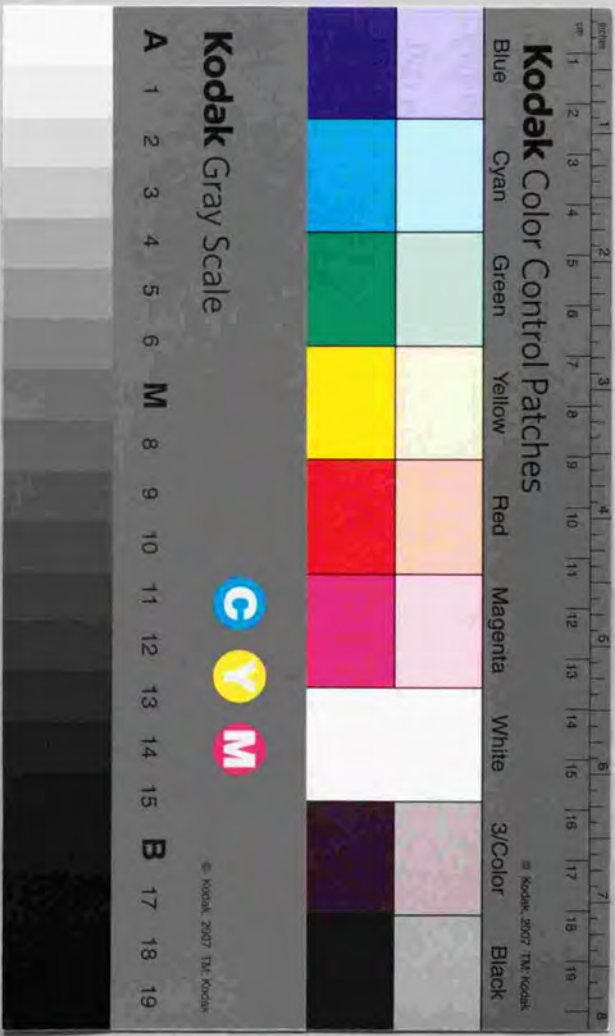
#### 国際会議

- (1) T. Aizawa, T. Ehara, N. Uratani, H. Tanaka, N. Kato, S. Nakase, N. Aruga and  
T. Matsuda: "A Machine Translation System for Foreign News in Satellite Broadcasting"  
COLING90, pp.308-310(1990)
- (2) N. Kato, N. Uratani and T. Aizawa:  
"Processing Proper Nouns in Machine Translation for English News"  
International Conf. on CICL, pp.431-439(1991)
- (3) J. Hosaka, T. Takezawa and N. Uratani:  
"Analyzing Postposition Drops in Spoken Japanese"  
ICSLP92, pp.1251-1254(1992)
- (4) T. Tashiro, N. Uratani and T. Morimoto:  
"Restructuring Tagged Corpora with Morpheme Adjustment Rules"  
COLONG94, pp.569-573(1994)
- (5) O. Furuse, Y. Sobashima, T. Takezawa and N. Uratani:  
"Bilingual Corpus for Speech Translation"  
AAAI-94 Workshop, (Integration of NL and SP), pp84-90, (1994)
- (6) T. Morimoto, N. Uratani, T. Takezawa, O. Furuse, Y. Sobashima, H. Iida,  
A. Nakamura, Y. Sagisaka, N. Higuchi and Y. Yamazaki:  
"a Speech and Language DataBase for Speech Translation Research"  
ICSLP94, pp.1791-1794(1994)

その他

- (1) 匂坂, 浦谷: "ATR 音声・言語データベース"  
日本音響学会誌, Vol.48, No.12, pp.878-882(1992)
- (2) 浦谷: "自動翻訳電話の現状"  
情報通信ジャーナル, 1993年3月号, pp.42-45(1993)
- (3) 浦谷, 江原: "音声処理と翻訳"  
テレビジョン学会誌, Vol.47, No.12, pp.1588-1591(1993)
- (4) "自動翻訳電話", 共著 (ATR 国際電気通信基礎技術研究所編)  
オーム社(1994)





**Kodak Color Control Patches**

Blue Cyan Green Yellow Red Magenta White 3/Color Black

**Kodak Gray Scale**

A 1 2 3 4 5 6 M 8 9 10 11 12 13 14 15 B 17 18 19



© Kodak, 2007 TM, Kodak