

LZW 圧縮テキストに対する高速文字列照合アルゴリズム

4 C-7

喜田 拓也 竹田 正幸 篠原 歩
九州大学大学院システム情報科学研究科

1. はじめに

圧縮されたテキストに対し文字列照合を行なう場合、一時的に展開したのちに、一般的な文字列照合アルゴリズムを用いるという方法が考えられる。この方法は単純ではあるが、圧縮テキストとその展開後のテキストの両方を走査しなければならないので時間がかかる。Amirら²⁾は、LZW符号により圧縮されたテキストを展開することなしに、文字列照合を行なうアルゴリズムを提案した。また著者らは、彼らのアルゴリズムをもとに、Aho-Corasickの文字列照合機械を利用することで、複数パターンすべての出現位置を出力するアルゴリズムを開発した⁴⁾。本稿では、Abrahamson¹⁾の文字列照合アルゴリズムを利用し、さらなる高速化と高機能化を図ったLZW圧縮テキストに対する文字列照合アルゴリズムを提案する。提案アルゴリズムは $O(m)$ 時間の前処理を必要とし、 $O(N+r)$ 時間でテキストを走査する。ここで m, N, r はそれぞれ、パターン長、圧縮テキスト長、パターンの出現回数である。

2. 一般化された文字列照合問題

Σ を文字の有限集合とする。文字列 ξ を構成する文字の数を ξ の長さといひ $|\xi|$ で表す。文字列 $\xi = a_1 a_2 \dots a_n$ の i 番目の文字 a_i を $\xi[i]$ で表し、 ξ の部分列 $a_i a_{i+1} \dots a_j$ を $\xi[i:j]$ で表す。

Abrahamson¹⁾は、次のように一般化された文字列照合問題を提案した。 $\Delta = \{X \subseteq \Sigma \mid \emptyset \neq X\}$ としたとき、

パターン $P = X_1 \dots X_m$ ($X_i \in \Delta$) と
 テキスト $T = T[1:n]$
 が与えられたとき、 $T[j:j+m-1] \in P$ となるような j をすべて求めよ。

このような j が一つでも存在するとき、テキスト T 中にパターン P が出現するといひ、この j をパターンの出現位置と呼ぶ。また、Abrahamsonは同時にこの問題を解く以下のアルゴリズムを示した。整数 p と整数の集合 S に対して $S \oplus p = \{i+p \mid i \in S\}$ とおく。 $1 \leq k \leq n$ に対して

$R_k = \{1 \leq i \leq m \mid X_1 \dots X_i \ni T[k-i+1:k]\}$
 と定義する。また、任意の $a \in \Sigma$ について、
 $M(a) = \{1 \leq i \leq m \mid X_i \ni a\}$

とおく。さらに関数 $f(S, a)$ を $f(S, a) = ((S \oplus 1) \cup \{1\}) \cap M(a)$ と定義しておく、このとき各 R_k は

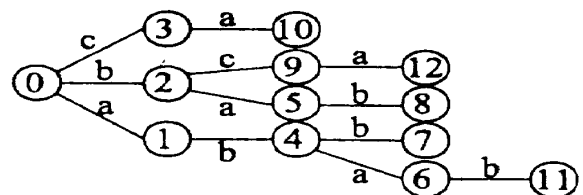
- 1) $R_0 = \emptyset$
- 2) $R_{k+1} = f(R_k, T[k+1])$ ($k \geq 0$)

によって求めることができる。アルゴリズムはテキストの文字を先頭から1文字ずつ読みながら、順次 R_k を計算していく。もし $R_k \ni m$ ならば $T[k-m+1:k] \in P$ である。すなわち $k-m+1$ がパターンの出現位置である。整数の集合を m ビットのビットベクトルで表現すると、関数 f はビットシフトと And, Or の論理演算によって定数時間で高速に求めることができる。またこのアルゴリズムには数多くの変種が開発されており、ミスマッチや複数パターンに対応したもの³⁾や、近似文字列照合が可能なもの⁵⁾が提案されている。

3. LZW 圧縮法

LZW圧縮されたテキストは整数の列になっている。各整数 j は辞書木とよばれる木構造の j 番目のノードに対応している。辞書木の各ノードは、木のルートからそのノードまでの枝にふられている文字の接続を表している。この辞書木のノードが示す文字列の集合を D と表し、辞書と呼ぶ。辞書木は、テキストを圧縮する際に適応的に構築される。また圧縮テキストを展開する際にも同じ辞書木を復元することが可能なので、辞書の情報を圧縮テキストに埋め込んでおく必要がない。図1は $\Sigma = \{a, b, c\}$, $T = abababbabcababcabab$ のときの辞書木を示している。例えば、この図の圧縮テキストにおける '9' は、文字列 bc を表している。よって以降では、文字列 u とそれを示す整数とを同一視する。

LZW圧縮テキストを展開するには、もとのテキストの長さに比例した計算時間が必要である。しかし、もとのテキストが必要ではなく、辞書木だけを構築するだけならば、その構築時間は $O(N)$ 時間ですむ(図2)。ここで N は圧縮テキストの長さである。



$T = abababbabcababcabab$
 $Z = 1, 2, 4, 4, 5, 2, 3, 6, 9, 11$

図1 辞書木

Input. LZW 圧縮テキスト $Z = u_1 u_2 \dots u_N$.

Output. 辞書木 D .

```

begin
  D := Σ;
  for i := 1 to N - 1 do begin
    if  $u_{i+1} \leq |D|$  then
      a を  $u_{i+1}$  の先頭の文字とする.
    else
      a を  $u_i$  の先頭の文字とする;
      D := D ∪ { $u_i \cdot a$ }
    end
  end
end.

```

図 2 圧縮テキストからの辞書木構築アルゴリズム.

4. 提案アルゴリズム

我々の目的は Abrahamson のアルゴリズムを, LZW 圧縮テキストに対して使用できるように拡張することである. LZW 圧縮テキストの 1 文字は辞書木の番号を示す整数であり, 辞書 D の文字列に対応することは既に述べた. まず, 任意の $v \in D$ について,

$$M^+(v) = \{1 \leq i \leq m \mid v \in \mathcal{P}[i - |v| + 1 : i]\}$$

と定義する. ただしこの式では $\mathcal{P}[i] = \Sigma$ ($i \leq 0$) としておく. この $M^+(v)$ を計算するために次の関係を用いる. 任意の $v \in D$, $a \in \Sigma$ について,

- 1) $M^+(a) = M(a)$,
- 2) $M^+(v \cdot a) = f(M^+(v), a)$

である. すなわち新たにノード u が辞書木に付加される時には, u の親ノード v の $M^+(v)$ から $O(1)$ 時間で $M^+(u)$ を計算することができる.

いま圧縮テキストが $Z = u_1 u_2 \dots u_N$ であるとする. このとき, $1 \leq \ell \leq N$ について,

- 1) $R_0 = \emptyset$,
- 2) $R_{k+|u_\ell|} = ((R_k \oplus |u_\ell|) \cup \{1, 2, \dots, |u_\ell|\}) \cap M^+(u_\ell)$

より R_k をとびとびに求めることができる. しかしながら, このようにして求めた R_k からはパターンの出現を検知できない場合がある. そこで次のような出力関数を用意する.

$$Output(R_k, u_\ell) = \{1 \leq i \leq |u_\ell| \mid R_{k+i} \ni m\}$$

するとこのとき, $\ell' = \sum_{j=1}^{\ell-1} |u_j|$ とおくと, パターンの出現位置は $\ell' + p - m + 1$ ($p \in Output(R_k, u_\ell)$) である.

任意の $u \in D$ について, $U(u) = \{1 \leq i \leq |u| \mid u[1:i] \in \mathcal{P}[m-i+1:m]\}$ とする. また \bar{u} を $U(u)$ の要素で一番大きな整数とし, 整数の集合 S に対して $m \ominus S = \{m-i \mid i \in S\}$ とする. このとき, 次の補題が成り立つ.

補題. $A(u) = \{\bar{u} \leq i \leq |u| \mid u[i-m+1:i] \in \mathcal{P}\}$ とおくと, このとき $Output(R_k, u_\ell)$ は

$$Output(R_k, u_\ell) = ((m \ominus R_k) \cap U(u_\ell)) \cup A(u_\ell)$$

と書くことができる. □

したがって, LZW 圧縮テキストに対する文字列照合アルゴリズムは図 3 のように要約できる.

Input. LZW 圧縮テキスト $u_1 u_2 \dots u_N$.

Output. すべてのパターンの出現位置.

```

begin
  ℓ' := 0, k := 0, R := ∅;
  for ℓ := 1 to N do begin
    UpdateDictionaryTrie( $u_\ell$ );
    for each  $p \in Output(R, u_\ell)$  do
      ℓ' + p - m + 1 にパターンが出現したと報告;
      R = ((R + |u_ℓ|) ∪ {1, 2, ..., |u_ℓ|}) ∩ M^+( $u_\ell$ );
      ℓ' := ℓ' + |u_ℓ|
    end
  end
end.

procedure UpdateDictionaryTrie( $u$ )
begin
  D に新しいノード  $u = v \cdot a$  を加える;
  M^+( $u$ ) = f(M^+( $v$ ), a);
  if M^+( $u$ ) ∋ m then
    if |u| ≤ m then U( $u$ ) := U( $v$ ) ∪ {|u|}
    else A( $u$ ) := A( $v$ ) ∪ {|u|}
    else U( $u$ ) := U( $v$ ), A( $u$ ) := A( $v$ )
  end;
end;

```

図 3 文字列照合アルゴリズム.

5. 実験

提案アルゴリズムがどの程度高速なのかを知るために, アルゴリズムを実際に実装し比較実験を行った. 図 4 がその実験結果である.

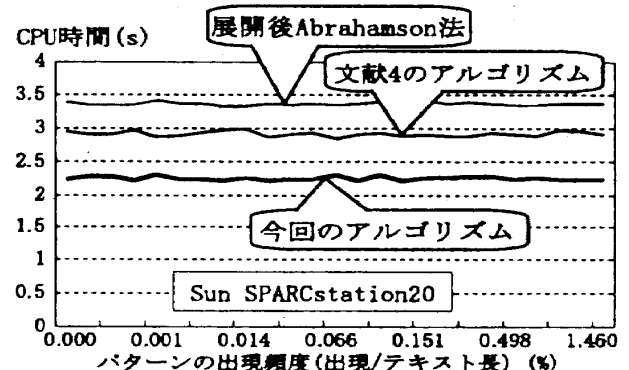


図 4 実験結果 (テキスト約 6.5Mb/圧縮時約 3.4Mb).

参考文献

- 1) K. Abrahamson. Generalized string matching. *SIAM J. Comput.*, 16(6):1039-1051, December 1987.
- 2) A. Amir, G. Benson, and M. Farach. Let sleeping files lie: Pattern matching in Z-compressed files. *Computer and System Sciences*, 52(23):299-307, 1996.
- 3) R. Baeza-Yates and G. H. Gonnet. A new approach to text searching. *Communications of the ACM*, 35(10):74-82, October 1992.
- 4) T. Kida, M. Takeda, A. Shinohara, M. Miyazaki, and S. Arikawa. Multiple pattern matching in LZW compressed text. In J. A. Atorer and M. Cohn eds., *Proceedings of Data Compression Conference '98*, pp. 103-112. IEEE Computer Society, March 1998.
- 5) S. Wu and U. Manber. Fast text searching allowing errors. *Communications of the ACM*, 35(10):83-91, October 1992.