



## 1. まえがき

自然言語処理システムは、辞書情報を利用した解析が必要不可欠であり、辞書検索効率がシステム全体の能率に大きな影響を与える<sup>(9, 20, 27)</sup>。特に、分かち書きされていないストリング処理(かな漢字変換, 日本語形態素解析)では、見出し(キー)の切り出しに有効なキー検索手法を利用する必要がある。

キー検索法<sup>(6, 21, 26)</sup>は適用分野により、キーの更新が頻繁に行なわれる動的キー集合を対象とする場合(動的検索法と呼ぶ)とこの更新が全く行なわれない静的キー集合を対象とする場合(静的検索法と呼ぶ)とに分類できる。動的検索法ではキーの更新を考慮するために、空間の使用効率あるいは検索効率が低下する場合があるが、静的検索法ではこれらの効率を極めて高い水準におくことが可能となる。動的検索法には、分散記憶法であるハッシュ法; 探索木法である2進木法, B\*木法; デジタル検索法(digital search, 広義には探索木法と考えられる)であるトライ(trie)法<sup>(2, 9, 9, 11, 13, 16, 17, 19, 21-26)</sup>, トライハッシュ法, 拡張ハッシュ法<sup>(15, 19)</sup>等がある。また、静的検索方法は、完全ハッシュ法<sup>(12, 14)</sup>が代表的であり、その外トライに対する静的データ構造<sup>(22)</sup>を提案したものがある。しかし、この分類は極端であって、現実の応用分野では、自然言語処理における辞書の様にキー集合の大部分が予め分かっている、後で適宜キーを追加して拡充する場合も非常に多い。この様なキー集合に対する検索法を準静的検索法<sup>(7)</sup>と呼ぶ。

従って、本報告では、ダブル配列法<sup>(7-10)</sup>を使ったトライ法を準静的検索法として確立するために、ダブル配列の拡張法を提案し、拡張されたダブル配列に対するキーの検索, 追加, 削除アルゴリズムを提案する。そして、これらアルゴリズムに対する最悪の計算時間を決定し、種々のキー集合に対する実験結果から本手法の有効性を実証する。

## 2. トライ検索マシンとダブル配列

### 2.1 トライ検索マシン

キー集合Kに対するトライ検索マシンMをストリングパターンマッチングマシン<sup>(1, 3, 4, 5)</sup>の記法で定義する。

$$M = (S, l, g, s_1, U)$$

で表す。ここで、S, lはそれぞれ状態, 入力記号の有限集合;  $s_1$ は初期状態(番号1で表す); gは状態遷移関数でgoto関数と呼ばれる。U( $\subseteq S$ )は出力状態の有限集合を表す。以後、状態は整数値の番号を対応させる。関数gは $S \times l$ から $S \cup \{fail\}$ への写像を定義し、状態sからtへの遷移が記号aに対して定義されていれば、 $g(s,$

$a)=t$ と書き、そうでなければ $g(s,a)=fail$ と書く。以後、状態sに入る遷移の数をindegree(s), そこから出る遷移の数をoutdegree(s)で表す。そして、 $outdegree(s)=0$ なる状態sを最終状態と呼び、その集合をFで表す。

また、goto関数gは $S \times l^*$ から $S \cup \{fail\}$ への写像に拡張して記述される。ここで、 $l^*$ はl上のすべてのストリングの集合を表す。以後、 $a, b, c, d \in l$ ;  $x, y, z \in l^*$ を使用し、空記号列(empty string)を $\epsilon$ で表す。以後、キーの最後尾にのみ付加される特別な記号#を定義し、最終状態と出力状態を一致させ、次の定義を行う。

[定義1] マシンMに対して、次の状態を定義する。

(1)  $g(s,a)=t, outdegree(s) \geq 2$ なる状態sから到達可能な最終状態までの遷移列上に $outdegree(r) \geq 2$ なる状態rが存在しないならば、状態sをセバレート状態と呼び、その集合を $S_p$ で表す。

(2) 初期状態からセバレート状態までの遷移列上に存在する状態(初期状態は含み、セバレート状態は含まない)をマルチ状態と呼び、その集合を $S_M$ で表す。

(3) セバレート状態から最終状態までの遷移列上に存在する状態(セバレート状態は含み、最終状態は含まない)をシングル状態と呼ぶ。

[定義2]  $g(s,y)=t, s \in S_p, t \in F$ なるストリングyをセバレート状態sに対するシングルストリングと呼び、STR[s]で表す。

[例1] キー集合 $K'=\{bac\#, bc\#, ba\#, bab\#$ に対するgoto関数を図1に示す。マルチ状態は1,3,7; シングル状態は4,5,6,9であって、セバレート状態と対応するシングルストリングは次のようになる。

STR[4]=#, STR[5]=#, STR[6]= $\epsilon$ , STR[7]=#.

### 2.2 ダブル配列

ダブル配列<sup>(6-10)</sup>はBASEとCHECKの二つの一次元配列を使用し、 $g(s,a)=t$ なる遷移を次の様に定義する。

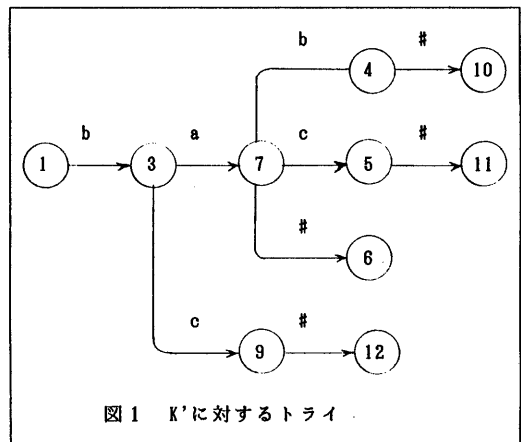


図1 K'に対するトライ

$$t = \text{BASE}[s] + a; \text{CHECK}[t] = s$$

ここで、記号 $a$ は1以上の内部表現値(numerical value)として取り扱われている。ダブル配列法は状態番号 $s$ と入力記号 $a$ の内部表現値との関係を利用して、遷移を確認すべきCHECKのインデックス(しかもこれが次に進むべき状態番号にも一致する)を巧妙に決定する(図2参照)。

本手法では、 $S_M \times (S_M \cup S_P)$ 上の遷移をダブル配列に格納し、セパレート状態以後の遷移は配列TAILにシングルストリングとして格納する。この拡張方法は、次の効果がある。

- 1) ダブル配列に格納する遷移数が減少するので、配列要素のバイト数を少なくできる。また、ダブル配列では一つの遷移に二つの要素を必要とするが、TAIL上では1記号のバイト数だけに記憶量を軽減できる。
- 2) シングル状態からの遷移がストリングのマッチングだけで確認できるので、検索を高速化できる。

しかし、この拡張に対しては処理中の状態がセパレート状態であるか否かを判定し、しかもTAIL中のシングルストリングの位置が決定できる効率的手法が必要不可欠となる。そこで、本報告では、次の条件を満足するダブル配列を定義する。

【定義3】 次が成立するならば、ダブル配列は条件Aを満足すると言う。

(A-1)  $g(s, a) = t; s \in S_M; t \in S_M \cup S_P$ ならばそのときに限って次が成立する。

$$\text{BASE}[s] + a = t, \text{CHECK}[t] = s.$$

(A-2)  $s \in S_P$ なる状態番号 $s$ に限って次が成立する。

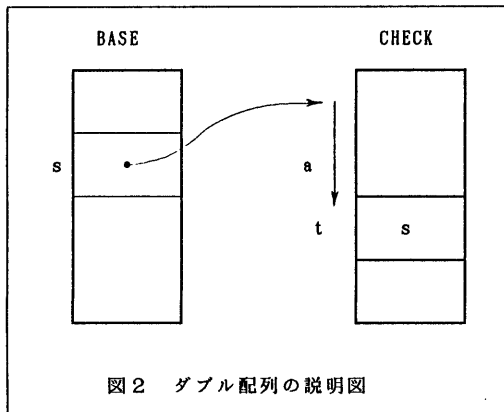
$$\text{BASE}[s] < 0.$$

(A-3)  $\text{STR}[s] = b_1 b_2 \dots b_m (1 \leq m)$ ;  $s \in S_P$ なるセパレート状態 $s$ に対して、次が成立する。

$$p = -\text{BASE}[s];$$

$\text{TAIL}[p] = b_1, \text{TAIL}[p+1] = b_2, \dots, \text{TAIL}[p+m-1] = b_m.$

条件(A-2)と(A-3)は $\text{BASE}[s]$ が負ならば $s$ はセパレート状態であって、しかも対応するシングルストリング



がTAILの位置 $-\text{BASE}[s]$ からアクセスできることを意味する。また、ダブル配列のインデックスの最大値を越える場合の判定を行なうために、次を定義する。ただし、ダブル配列の未使用要素の値0とする。

【定義4】 次が成立するならばダブル配列は条件Bを満足すると言う。

初期状態番号1に対する $\text{CHECK}[1]$ の値は、ダブル配列の $\text{CHECK}[q] \neq 0$ なる最大のインデックス $q$ と一致する。

### 3. ダブル配列によるキー検索

拡張されたダブル配列を使用した検索アルゴリズム1を図3に示す。但し、次の変数と関数を使用する。  
 $S\_TEMP$ : TAILから取り出されたストリングを保持する。  
 $\text{FETCH\_STR}(p)$ :  $\text{TAIL}[p]$ から $\text{TAIL}[p+k] = \#, 0 \leq k$ なる記号( $\#$ は含む)までのストリングを返す。

$\text{STR\_CMP}(y, z)$ : ストリング $y$ と $z$ が一致すれば-1を、一致しなければ部分マッチした接頭辞の長さを返す。

アルゴリズム1の行(1-1)は、 $t$ がCHECKの最大のインデックスを越えるかあるいは遷移が未定義の場合、FALSEを返す。また、repeat文は $\text{BASE}[s]$ が負になれば終了し、行(1-2)はシングルストリング $\text{STR}[s]$ を $S\_TEMP$ に代入する(ただし、 $h=n+1$ ならばマッチングが成功し、TRUEを返す)。そして、行(1-3)で $S\_TEMP$ と残りの入力ストリングとの比較により $x$ が $K$ に含まれるか否かを決定する。

【例2】 キー集合 $K'$ に対するダブル配列(次章で構成される)を図4に示す。ただし、記号 $a, b, c, \#$ の内部コードはそれぞれ1, 2, 3, 4とし、TAILの記号'?'は次章で説明される未使用記号を表す。キー $bc\#$ の検索例を示す。  
 $t = \text{BASE}[1] + b = 3$

**Algorithm 1.** Retrieval Algorithm.

**Input:** A string  $x = a_1 a_2 \dots a_n a_{n+1} a_{n+1} = \#;$   
 and the machine  $M$  with the double-array and TAIL for  $K$ .

**Output:** If  $x \in K$ , then the output is TRUE, otherwise FALSE.

**Method:**

```

begin
  s:=1; h:=0;
  repeat
    h:=h+1; t:=BASE[s]+ah;
  (1-1) if (t>CHECK[1]) or (CHECK[t]≠s)
    then return(FALSE) else s:=t
  until BASE[s]<0;
  (1-2) if ah=# then return(TRUE)
  else S_TEMP:=FETCH_STR(-BASE[s]);
  (1-3) if STR_CMP(ah+1... anan+1, S_TEMP)=-1
    then return(TRUE) else return(FALSE)
end.
```

図3 ダブル配列による検索アルゴリズム

	1	2	3	4	5	6	7	8	9
BASE	1	0	6	-8	-1	-7	2	0	-5
CHECK	9	0	1	7	7	7	3	0	3
	1	2	3	4	5	6	7	8	9
TAIL	#	\$	?	?	#	\$	\$	#	\$

図4 K'に対するダブル配列

- (1-1):  $t=3 \leq \text{CHECK}[1]=9$  且つ  $\text{CHECK}[3]=1$   
 $\text{BASE}[s]=\text{BASE}[3]=6 > 0$ ,  $t=\text{BASE}[3]+c=9$
- (1-1):  $t=9 \leq \text{CHECK}[1]=9$  且つ  $\text{CHECK}[9]=3$   
 $\text{BASE}[s]=\text{BASE}[9]=-5 < 0$
- (1-2):  $h=2 \neq n+1=3$ となるので、関数FETCH\_STRにより  
 $S\_TEMP = \#$ がセットされる。
- (1-3): 関数STR\_CMPで残りの入力ストリング $\#$ と  
 $S\_TEMP$ が比較され、 $bc\# \in K'$ が判定される。

#### 4. ダブル配列の更新

##### 4.1 キーの追加

ダブル配列の初期状況は、 $\text{CHECK}[1]$ と $\text{BASE}[1]$ が共に1であって、最大インデックス $\text{CHECK}[1]$ を越える要素の値は0とする。このとき、追加アルゴリズム2はアルゴリズム1の行(1-1, 1-3)のreturn(FALSE)を次のように変更することで得られる。

a) 行(1-1)のreturn(FALSE)の変更:

```
begin
  A_INSERT(s, ahah+1 ... anan+1);
  return(FALSE)
```

end

b) 行(1-3)のreturn(FALSE)の変更:

```
begin
```

/\* S\_TEMPを $a_{h+1}a_{h+2} \dots a_{h+r}b_1 \dots b_m$ , 残りの入力ストリングを $a_{h+1}a_{h+2} \dots a_{h+r}a_{h+r+1} \dots a_n a_{n+1}$ として、 $b_1 \neq a_{h+r+1}, 0 \leq r \leq n-h+1, 0 \leq m$ と仮定する\*/

```
B_INSERT(s, ah+1 ... ah+r,
          ah+r+1 ... anan+1, b1 ... bm);
  return(FALSE)
```

end;

アルゴリズム2で使用する主要な手続きと関数を図5~7に示す。但し、次の変数と関数を使用する。

SET\_LIST(s):  $g(s, a) \neq \text{fail}$ なる記号aを要素とする集合を返す。

LIST, LIST1, LIST2: {#} ∪ Iの部分集合。

N(LIST): LISTの要素数を返す。

POS: TAILの長さに1を加えた値を保持するグローバル変数であって、初期値は1である。

current: MODIFYで変更された状態番号sに対する新しい番号を保持する(詳細は以下で説明される)。

X\_CHECK(LIST):  $c \in \text{LIST}$ なる全てのc(nullは除く)が $\text{CHECK}[q+c]=0$ を満足する最小のインデックスqを返す。

SET\_STR(p, y): ストリングyをTAILのインデックスpの位置から連鎖し、pがPOSと等しければPOSにyの長さを加えた値を、そうでなければPOSの値を返す。

STR\_LEN(z): ストリングzの長さを返す。

```
procedure A_INSERT(s, b b ... b );
                1 2      m
begin
(a-1) t:=BASE[s]+b ;
(a-2) if CHECK[t]≠0 then
      begin
(a-3) LIST1:=SET_LIST(s);
(a-4) LIST2:=SET_LIST(CHECK[t]);
(a-5) if (N(LIST1)+1 < N(LIST2)) then
(a-6) s:=MODIFY(s, s, b , LIST1)
      else
(a-7) s:=MODIFY(s, CHECK[t], null, LIST2);
      end;
(a-8) INS_STR(s, b b ... b , POS);
                1 2      m
end;
```

図5 手続きA\_INSERT

```
function MODIFY(current, s, a, LIST);
begin
(m-1) old_base:=BASE[s];
(m-2) BASE[s]:=X_CHECK(LIST ∪ {a});
      for each c in LIST do
      begin
(m-3) t:=old_base+c; t':=BASE[s]+c;
(m-4) CHECK[t']:=s; BASE[t']:=BASE[t];
(m-5) if BASE[t]>0 then
      begin
(m-6) for each q such that CHECK[q]=t do
(m-7) CHECK[q]:=t';
(m-8) if (t=current) current:=t';
      end;
(m-9) BASE[t]:=0; CHECK[t]:=0
      end
      return(current)
end;
```

```
procedure INS_STR(s, e e ... e , d_pos);
                1 2      q
begin
(s-1) t:=BASE[s]+e ;
(s-2) CHECK[t]:=s; BASE[t]:=-d_pos;
(s-3) POS:=SET_STR(d_pos, e e ... e $);
                2 3      q
end;
```

図6 関数MODIFYと手続きINS\_STR

```

procedure B_INSERT(s, x, y, z);
/*x=b1b2... bk; y=c1c2... ch;
z=d1d2... du*/
begin
(b-1) old_pos := -BASE[s];
    for i:=0 to STR_LEN(b1b2... bk) do
    begin
(b-2) BASE[s]:=X_CHECK({bi});
(b-3) CHECK[BASE[s]+bi]:=s;
(b-4) s:=BASE[s]+bi
    end;
(b-5) BASE[s]:=X_CHECK({c1,d1});
(b-6) INS_STR(s,d1d2... du, old_pos);
(b-7) INS_STR(s,c1c2... ch, POS)
    end;

```

図7 手続き B\_INSERT.

手続き A\_INSERT は、ダブル配列上でキーがミスマッチしたときにキーを追加し、もし BASE[s] の変更が必要ならば、次の優先順位で MODIFY を呼ぶ。

$N(LIST1)+1$  (状態 s からの遷移数と新しい遷移  $g(s, b_1)$  の和) と  $N(LIST2)$  (CHECK[t] から出る遷移数) を比較して、少ない方の状態番号の BASE の値を変更する (詳細は定理 1 の証明と例 3 を参照)。

手続き INS\_STR は、セバレート状態に至る遷移を定義し、シングルストリングを TAIL に格納する。また、手続き B\_INSERT は、キーが TAIL 上でミスマッチしたときにキーを追加する。第 2 引き数は、残りの入力ストリング xy と S\_TEMP(xz とする) で部分マッチした接頭辞 x; 第 3, 4 引き数は xy, xz からこの接頭辞 x を除いたストリング y, z をそれぞれ表す。

ダブル配列の大きさ CHECK[1] は関数 X\_CHECK で変更されるものとし、以後条件 B の議論は省略する。

[定理 1] アルゴリズム 2 は、 $x \in K$  なる入力 x をダブル配列に追加し、得られた配列は条件 A を満足する。

(証明) 各手続きの正当性を証明する。

### 1) 手続き A\_INSERT.

行 (a-2) の条件式が偽の場合、行 (a-8) で呼ばれる INS\_STR の行 (s-1, s-2) は遷移  $g(s, b_1)$  をダブル配列に定義し、行 (s-3) はシングルストリングを TAIL に格納する。しかも、t は必ずセバレート状態番号となるので、条件 A は明らかに保存される。また、行 (a-2) が真で、行 (a-5) の比較により行 (a-6) の MODIFY が呼ばれたとき、条件 (A-1) の成立は次の 3 点について証明すればよい。

#### a) 新しい BASE[s] による遷移の定義.

行 (m-2) の X\_CHECK が決定する新しい BASE[s] は、 $b \in LIST \cup \{a\}$  なる全ての b に対して  $CHECK[BASE[s]+b]=0$  を満足し、行 (m-4) は  $g(s, c)$ 、 $c \in LIST$  なる遷移を定義する。また、行 (m-8) は  $current(=s)$  を変更しないので、

行 (a-8) で呼ばれる INS\_STR の行 (s-1, s-2) は新しい遷移  $g(s, b_1)$  を正しくダブル配列上に定義する。

#### b) 他の遷移への影響.

BASE[s] の変更に伴い行 (m-3) の番号 t は新しい番号 t' 変化するが、行 (m-4) が  $BASE[t]$  を  $BASE[t']$  にコピーし、行 (m-7) が  $CHECK[BASE[t]+b]=t$  なる CHECK の要素を t' に変更するので、矛盾は生じない。

#### c) 古い遷移の除去

行 (m-9) により、明かである。更に、条件 (A-2, A-3) の成立は、行 (a-8) で呼ばれる INS\_STR の行 (s-2 ~ s-3) より明かである。また、行 (a-7) で呼ばれる MODIFY は、BASE[s] の代わりに  $BASE[CHECK[t]]$  の値を変更するので、番号 s が変更される可能性がある。これに対して、第 1 引き数には番号 s が current として送られ、s が変更された場合は行 (m-8) が current を新しい番号 t' に変更するので、矛盾は生じない。

### 2) 手続き B\_INSERT.

行 (b-2 ~ b-4) は  $b_1 b_2 \dots b_k$  に対する遷移を定義し、行 (b-5) 及び行 (b-6, b-7) で呼ばれる INS\_STR の行 (s-2) は  $c_1$  と  $d_1$  に対する遷移を定義する。次に、行 (b-6) で呼ばれる INS\_STR の行 (s-3) は  $d_2 \dots d_m$  を TAIL[old\_pos] (old\_pos は行 (b-1) で保持する) からオーバーライトする。また、行 (b-7) で呼ばれる INS\_STR の行 (s-3) が  $c_2 \dots c_n$  を TAIL[POS] から新たに格納する。しかも、これらのシングルストリングの格納において、行 (s-2) が  $BASE[t]$  をそれぞれ -old\_pos と -POS に設定する。従って、条件 A の成立は明かである。

故に、定理は成立する。 (証明終)。

[例 3] キー bac#, bc#, ba#, bab# の追加過程を図 8 に示す。ただし、最終結果は既に図 4 に示されている。また、行 (s-3) のオーバーライトで生じた TAIL の未使用要素を '?' で表す。以下にキー bc#, bab# の処理を行番号と共に簡単に記述する。

bc# では、B\_INSERT(3, ε, c#, ac#) が実行される。

(b-1): old\_pos = -BASE[3] = 1

(b-5): BASE[3] = X\_CHECK({c, a}) = 1

(b-6): INS\_STR(3, ac#, 1)

(s-3): POS = SET\_STR(1, c#\$) = 5

(b-7): INS\_STR(3, c#, 5)

(s-3): POS = SET\_STR(5, #) = 7

bab# では A\_INSERT(2, b#) が実行される。

(a-1, a-2): t = BASE[2] + b = 4, CHECK[4] = 3 ≠ 0

(a-3): LIST1 = SET\_LIST(2) = {c, #}

(a-4): LIST2 = SET\_LIST(3) = {a, c}

(a-5):  $N(LIST1)+1=3 > N(LIST2)=2$

(a-7): MODIFY(2, 3, null, {a, c})

(m-1): old\_base = BASE[3] = 1

(m-2): BASE[3] = X\_CHECK({a, c}) = 6

LIST の要素 a に対して、次が実行される。

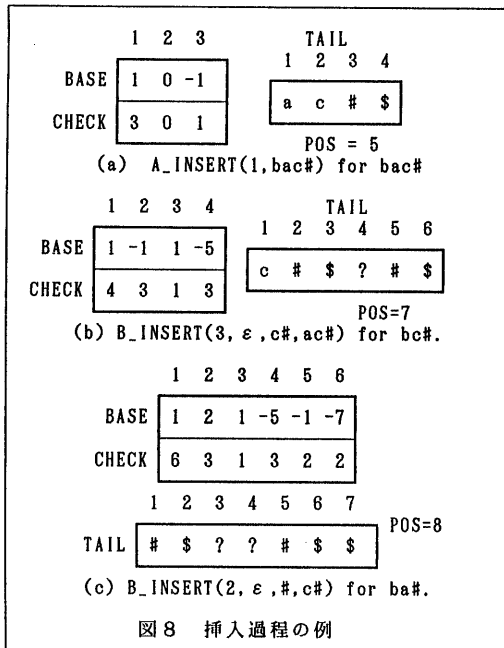


図8 挿入過程の例

- (m-3):  $t = \text{old\_base} + a = 2$ ,  $t' = \text{BASE}[3] + a = 7$
- (m-4):  $\text{CHECK}[7] = 3$ ,  $\text{BASE}[7] = 2$
- (m-5):  $\text{BASE}[2] = 2 > 0$
- (m-7):  $\text{CHECK}[5] = 7$ ,  $\text{CHECK}[6] = 7$
- (m-8):  $\text{current} = 7$  (現在の状態番号2の変更)
- (m-9):  $\text{BASE}[2] = 0$ ,  $\text{CHECK}[2] = 0$

LISTの要素cに対する処理は省略するが、MODIFYは状態番号7を返し、以下が実行される。

- (a-8):  $\text{INS\_STR}(7, d\#, 8)$
- (s-3):  $\text{POS} = \text{SET\_STR}(8, \#\$, 10)$

#### 4.2 キーの削除

削除アルゴリズムは、アルゴリズム1の行(1-2, 1-3)のreturn(TRUE)を次の様に変更することで得られる。行(1-2, 1-3)のreturn(TRUE)の変更：

```
begin
    BASE[s]:=0; CHECK[s]:=0;
    return(TRUE)
end
```

[定理2] アルゴリズム3は、ダブル配列から  $x \in K$  なるキーを削除し、得られたダブル配列は条件Aを満たす。

(証明) 負の状態番号は存在しないので、 $\text{CHECK}[s] = 0$ によりセパレート状態s, に入る  $g(s, a) = s$ , なる遷移はダブル配列上で明らかに未定義となる。しかも、セパレート状態s, はキーxに対して唯一であるので、定理の成立は明かである。 (証明終)。

## 5. 評価

キーの総数をn, ダブル配列に定義されるべき状態数 ( $S_M \cup S_P$ の要素数)をuとすると、 $u = c_1 n$ なる定数  $c_1$  は存在するが、 $m = c_2 u$ なる定数  $c_2$  はダブル配列の構成状況により一意に決定できないので、nとeによるmの厳密な評価は難しい。ただし、あえてnとeで評価するならばmは  $u \cdot e$  を越えることはないので、正確な近似ではないが  $O(m)$  は  $O(n \cdot e)$  と表せる。従って、理論的評価ではこのmを使用し、mには具体的評価を与える。なお、理論的評価はすべて最悪の場合とする。

### 5.1 理論的評価

以上、アルゴリズムの記述を簡単にするためにダブル配列の未使用要素の値は0として議論してきたが、X\_CHECKの計算時間を決定するために、次を定義する。

ダブル配列の未使用状態番号(未使用要素のインデックス)を昇順に  $r_1, r_2, \dots, r_m$  とするとき、

$$\begin{aligned} \text{CHECK}[r_i] &= -r_{i+1}; \quad 1 \leq i \leq m-1 \\ \text{CHECK}[r_m] &= -1 \end{aligned}$$

なるリンクを作る。ただし、 $r_i$  は変数G\_HEADで与える。このG\_HEADからのリンク(G-リンクと呼ぶ)によりX\_CHECKの計算時間は、 $O(m \cdot e)$  となる。

検索の計算時間は、キーの長さkに比例し  $O(k)$  となり、また削除の計算時間(検索時間は省略)は、未使用状態番号をG-リンクへ追加する操作に依存し、 $O(m)$  となる。また、追加の計算時間(検索時間は省略)は、B\_INSERTではなくA\_INSERTのMODIFYに依存する。まず、X\_CHECKは  $O(m \cdot e)$  で、二つのfor文のループ回数の最大値はeとなるので、追加の計算時間は  $O(m \cdot e + e^2)$  となる。

ダブル配列の大きさは  $O(n+m)$  となる。また、TAILの未使用記号 '?' は、キーを辞書順に追加する場合にはなくすることが可能である。即ち、B\_INSERTの行(b-6)で呼ばれたINS\_STRで再格納されるシングルストリングは先に格納されていたものより短くなり、常にTAILの最後尾でのみオーバーライトされるので、SET\_STRの返す値を短くなった分だけPOSから引いてやればよい。

### 5.2 具体的評価

提案されたダブル配列によるキー検索システムは、言語Cで記述されており、1) 設計ルーチン; 2) 検索、追加、削除ルーチン; 3) 記憶管理ルーチンから成っており、それぞれ約300行である。1)は入力記号の種類と符号化(英数、片仮名等); キー集合の分割基準(ハッシュ法、各キーの数文字分の接頭辞による部分デジタル検索木等で定義できる)を設定する; TAILの設計(レコード情報の取り扱い等)を行う。この分割は、ダ

ブル配列要素が2バイトを越える場合に実行して、記憶領域を小さくする効果がある。また、パソコンの様に主記憶が小さくダブル配列とTAILを補助記憶からアクセスする場合は分割数を多くして、効率的な転送を実現できる。2)は、提案された検索、追加、削除のルーチンである。3)は分割されたダブル配列とTAILの記憶管理(ダブル配列とTAILの各組の格納番地と大きさをもつ表の管理)を実現する。従って、分割を行った場合は、各キーに対応するダブル配列とTAILの組の決定手順とこの記憶管理が必要となる。このシステムは、Sun Sparc Station, NEC PC-9801の種々のコンピュータ上で稼働している。

表1に示す実験結果の $K_1, K_2, K_3, K_4, K_5$ はそれぞれPascalの指定語; COBOLの指定語; 世界の主要都市名地名; 英単語辞書(派生語約6万の中心語によるキー); 日本語辞書(約5万語表記の片仮名読みによるキー)であり、特に $K_5$ は約30種類の実験結果の中から最も条件の悪い結果を与えるものとして採用した。ただし、 $K_4$ と $K_5$ はダブル配列を五つに分割した結果であり、また記号'#'の内部コードは1に対応させてある。

表1のシングル状態数は全体の約20~70%( $k_5$ を除けば50~70%)となり、TAILの有効性が分かる。記憶量の情報は、ダブル配列(TAILを含む)、TAIL('\$'は含まない)、ソース記憶量(区切り記号を含むキー集合の大きさ)を表す。 $K_3$ 以上のキー数では共通の接頭辞の効果と極端に低い未使用率(ダブル配列中の未使用要素の占有率)でダブル配列とTAILの大きさがソース記憶量の約1.1~1.2倍にしかならないことが分かる。この未使用

率は、 $m$ が非常に小さい値になるという具体的評価を意味する。

ここで、ハッシュ法(衝突を鎖法で解決し、バケットの大きさは1とする)と2進木との記憶量の比較を行う。長さの異なるキーを検索するので、キーの代わりにキーファイルへのポインタ( $p$ で表し、4バイトとする)を使用し、ハッシュ表の大きさをキーの総数の4分の1、その要素を2バイトと仮定すれば、バケットはポインタ $p$ と鎖用のポインタ(2バイト)を; 2進木の各ノードはポインタ $p$ と二つのノードポインタ(2バイト)をもつ。この計算によりハッシュ法と2進木でキー検索を実現するための記憶量( $K_3$ 以上を対象にする)は、ソース記憶量のそれぞれ約1.7~2.2, 1.8~2.4倍となるので、本手法の空間的有効性が分かる。

追加と検索時間(VAXによる値)は、主記憶上にあるキー集合を辞書順に一括して追加し、その後同じキー集合を一括検索した場合の1件当たりの平均時間を表す。但し、 $K_1, K_2$ の検索時間は誤差を含む。この結果より、キー検索時間の高速性が分かる。また、追加時間も平均約35ms程度となり、十分高速であることが分かる。

追加の評価の $AVE-me, AVE-e, AVE-e', AVE-m$ は、それぞれX\_CHECKにおけるダブル配列要素の参照回数; MODIFYにおけるLISTの要素数; LISTの全要素に対して行( $m=5$ )が真となる回数;  $m$ の平均値を表し、()内それらの最大値を表す。本実験での $e$ の値は $K_5$ のみ63で、その他は96である。AVE-meより $m \cdot e$ の項は $K_5$ でも平均19.5(最大255)と決して大きい値ではない。また、AVE-e'がAVE-eよりも更に小さくなるので、 $e^2$ の項は $K_5$ の場合でも平均132.3(最大1,008)となる。従って、これらの最大値が同時に起こったとしてもその値は1,263となり、実用的な数値であることが分かる。また、AVE-mより $m$ の値は、追加が連続するときには常に小さい値となることが分かる。

## 6. むすび

本報告では、準静的検索法と呼ばれる基準を設定し、ダブル配列によるトライ検索でこれを実現した。そして、本手法が次の点を満足することを理論的且つ具体的に示した。

- 1) キーの数に関係しない高速な検索が可能である。
- 2) 空間の使用効率が高い(キー集合の大きさ程度)。
- 3) 実用的な時間内でキーの追加が可能であり、しかも1)と2)の特徴を低下させない。

また、ダブル配列は如何なるキー集合に対しても構成できるので、これらの特徴を持たない完全ハッシュ法<sup>(14)</sup>、線形ハッシュ法<sup>(27)</sup>、B木法<sup>(20)</sup>; トライ検索<sup>(11, 24)</sup>よりも広い適用範囲をもつキー検索法といえる。

表1 実験結果

各 値	$K_1$	$K_2$	$K_3$	$K_4$	$K_5$
各種カウント					
キーの数	35	310	1,480	23,976	32,344
キーの平均長	5.1	6.5	8.5	7.2	4.6
マルチ状態数	17	301	981	16,518	13,039
シングル状態数	74	637	5,801	35,765	11,092
全状態数	126	1,248	8,262	76,257	56,457
記憶量(K-bytes)					
ダブル配列	0.36	3.63	17.17	221	226
TAIL	0.11	0.95	7.28	60	44
ソース記憶量	0.18	2.33	14.08	196	183
未使用率(%)	17.5	8.9	0.48	0.23	0.41
時 間(ミ秒)					
追加	31.4	30.6	29.1	36.7	35.5
検索	2.9	0.97	0.41	0.37	0.35
追加の評価					
AVE-me	1.1(3)	4.7(31)	3.6(37)	3.1(44)	19.5(255)
AVE-e	2.0(3)	2.0(7)	1.8(13)	2.5(12)	3.5(21)
AVE-e'	0.26(1)	1.5(5)	1.8(8)	1.4(10)	2.1(16)
AVE-m	1.1(3)	3.1(24)	2.4(23)	2.1(32)	11.3(93)

本手法は、日英文章の解析システムに於ける辞書、高頻度英単語のフィルタリング等の約30種類のキー集合に応用されている。特に、ダブル配列は効率的デジタル検索を実現するので、日本語解析に於ける最長一致法と英語の語彙解析におけるスベルチェックと訂正の候補<sup>(24)</sup>等の実現にも有効なものとなった。

本論文では入力記号の符号化に付いて議論しなかったが、2進符号化を導入して効率的なトライハッシング<sup>(15,19)</sup>を実現することは、最悪の追加時間を線形にするためにも興味ある課題である。

謝辞：本研究の一部は、文部省科学研究費一般研究(C)、民間企業（株）ジャストシステムとの共同研究、（社）システム総合研究所委託研究、電気通信財団研究費などの援助を受けた。なお、本辞書検索手法は、多くの研究所に無料配布され、実用化されている。配布希望者は、770 徳島市南常三島町2-1, 徳島大学工学部知能情報工学科, Tel. 0886-23-2311 内4752, Fax. 0886-55-4424に連絡されたい。

#### 文 献

- (1) A.V.Aho et al., "Efficient string matching: An aid to bibliographic search", *Commun.ACM.*, 18, 6, pp.333-340 (1975).
- (2) A.V.Aho et al., *Data Structures and Algorithms*, Addison-Wesley, Reading Mass. Ch. 5 (1983).
- (3) A.V.Aho et al., *Compilers -Principles, Techniques, and Tools-*, Addison-Wesley, Reading Mass. Ch.3 (1986).
- (4) J.Aoe et al., "A method for improving string pattern matching machines", *IEEE Trans. Softw. Eng.*, SE-10, 1, pp.116-120 (1984)
- (5) J. Aoe, "An Efficient Implementation of Static String Pattern Matching Machines," *IEEE Trans. Softw. Eng.*, SE-15, 8, pp.1010-1016 (1989).
- (6) J. Aoe, *Computer Algorithms -Key Search Strategies-*, IEEE Computer Society Press. (1991).
- (7) J. Aoe, "An Efficient Digital Search Algorithm by Using A Double-Array Structure," *IEEE Trans. Softw. Eng.*, SE-15, 9, pp.1066-1077 (1989).
- (8) 青江, "ダブル配列による高速テイク検索アルゴリズム," *信学論(D)*, J71-D, 9, pp.1592-1600 (1988).
- (9) 青江, "自然言語辞書の検索—ダブル配列による高速テイク検索アルゴリズム—," *bit*, 21, 6, pp.776-784 (1989)
- (10) 青江, "ダブル配列による有限状態機械の効率的イ  
プリメンテーション," *信学論(D)*, J70-D, 4, pp.653-662 (1987).
- (11) A. W. Apple et al., "The world's fastest scrabble program," *Commun. ACM*, 31, 5, pp.572-578 (May, 1988).
- (12) F. Berman et al.: "Collections of functions for perfect hashing", *SIAM J. Comput.*, 15, 2, pp.604-618 (1986).
- (13) A. Blumer et al., "Complete inverted filters for efficient text retrieval and analysis," *Journal of ACM.*, 34, 3, pp.578-595, (Mar, 1987).
- (14) G. V. Cormack et al., "Practical perfect hashing", *The Comput. J.*, 28, 1, pp.54-58 (1985).
- (15) R. Fagin et al., "Extensible hashing-A fast access method for dynamic files," *ACM. Trans. Database Syst.*, 4, 3, pp.315-344 (1979).
- (16) E. Fredkin, "Trie memory," *Commun. ACM*, 3, 9, pp.490-500 (1960).
- (17) G.H.Gonnet et al., *Handbook of Algorithms and Data Structures*, Addison-Wesley, (1991).
- (18) R.F. Ilson(Editor): *Longman Dictionary of Contemporary English*, Longman Group Ltd. (1978).
- (19) W. Jonge et al., "Two access method using compact binary trees", *IEEE Trans. Softw. Eng.*, SE-13, 7, pp.799-809 (1987).
- (20) 日高外, "拡張B-Treeと日本語単語辞書への応用," *信学論(D)*, J87-D, 4, pp.399-404 (1984).
- (21) D.E.Knuth, *The Art of Computer Programming*, Vol.3, *Sorting and Searching*, Ch.6 (1973).
- (22) M. Kurt, "Compressed Tries", *Commun.ACM*, 19, 7, pp.409-415 (1976).
- (23) 中嶋外, "TRIE構造とガラスタックを用いた日本語形態素解析," *情処第39回全国大会*, IF-4, pp.589-590 (1989).
- (24) J.L.Peterson, *Computer Programs for Spelling Correction*, *Lecture Notes in Comput. Sci.*, Springer-Verlag, N. Y. (1980).
- (25) 田中, *自然言語解析の基礎*, 産業図書 (1989).
- (26) S.A.Thomas, *Data Structure Techniques*, Addison-Wesley, Reading Mass. Ch.3 (1980).
- (27) 横山外, "二次記憶上の大規模語彙を用いる自然言語処理システム," *情処論*, 29, 6, pp.570-580 (1988)
- (28) 鐘外, "Prefix-Closed B-Tree, 情処自然言語研, NL73-15, (1989).