

2005年12月作成

圧縮索引とその周辺

岡野原 大輔

hillbig@is.s.u-tokyo.ac.jp

東京大学

発表の流れ

- 背景

- 索引 / 全文索引 / 圧縮索引

- 基礎知識

1990年代

- Suffix Arrays, Burrows Wheeler Transform

- 圧縮全文索引

2000年

- Compressed SA, FM-index

- 最近の話題

2005年 ~

- Wavelet Tree, XWT (TreeのBWT)
- Succinctなデータ構造 (bit array, tree)



背景

大規模データの利用

- 非常に大きなテキストデータが手に入るようになった
 - MEDLINE (1100万アブストラクト 500GB)
 - Blog Watcher (1100 blog エントリー)
 - TREC2004 Terabyte Track (2500万文書 426GB)
 - Web Pages in Internet (~ PB)
 - Genome 配列 (> 800G 塩基対 in 2004)
- 「*We can obtain accurate information from very large inaccurate data (e.g. 50% is error) by central limit theorem*」 (Mehran Shaami@AAMT2005)
 - c.f. Googleの大規模なウェブデータを使った統計的機械翻訳システムは2005年のワークショップで最高精度。

問題設定

- **前もって**与えられたデータ T (長さ: n アルファベット集合:) とパターン P (長さ m)
- 次の二つの操作が答えられるか
 - $\text{occ}(P)$ パターン P のデータ中の出現回数
 - $\text{loc}(P)$ パターン P のデータ中の全ての出現位置
 - 他の多くの操作もこれらの基本操作の組み合わせに帰着される

用語

- 索引
 - パタンの出現した回数、場所を前もって求め保存したもの
- 全文索引
 - 全てのテキスト中に出現したパタンの出現した回数、場所を前もって求め保存したもの
 - 文字索引 (Suffix Arrays, Suffix Trees, N-gram など)
 - 任意の文字列を検索可
 - 単語索引 (転置ファイルなど)
 - 形態素解析等で抽出された単語の出現位置のみ記録
- 圧縮全文索引
 - 全文索引を圧縮したもの。復元操作は、前もって行われない
- 本発表では、圧縮全文索引の中で文字索引を扱う

なぜ索引？

- 大規模データを利用するには線形時間操作でさえコストが大きい
 - 例 データ1G中の“the”の出現場所を全て報告
 - $O(\log N)$ (suffix arrays) 0.005 msec
 - $O(N)$ (grep) 970 msec
- 索引を使うことで、パターンPの出現位置や回数を高速に求めることが可能

なぜ文字索引？

- 文字索引: 任意のパターンに対し答えられる
- 従来の転置ファイルは答えが漏れる場合がある
- 日本語等、分かち書きでない言語では特に問題。英語でもフレーズを探索できない
(統計的機械翻訳とかで問題)
- ゲノムなど単語が無い場合はさらに困難
- N-gram方式(N文字転置ファイル)が最近利用されている
 - この場合、辞書ヒット数が大きい場合、計算量は $O(N)$ に近づくため、大規模データでは使えない

なぜ圧縮全文索引？

- 索引では速度と作業領域量はトレードオフ関係
 - 極端な話、全ての答えを前もって求めおくと
計算量は $O(1)$ だが、作業領域量は非常に大きい
- 全文文字索引では、作業領域量が非常に大きい
 - 扱えるデータサイズに制限があった。
- このトレードオフはどっちを重視すればいいの？
作業領域量、速度の両面でほぼ最適なものが達成可能
 - 最新の結果: NH_k bitの作業領域量で $\text{occ}(P)$ を $O(m)$ time、
[Ferragina 2005]
 - H_k : 入力 T の k 次エントロピー。
 - 実データの H_5 の例 英文:0.23 DNA:0.24 XML:0.10

圧縮索引の応用例

- 情報抽出

- 固有表現抽出

$df_2(P)/df(P)$ [Church 2001]

- 生物情報

- ゲノム解析 (ゲノム比較、ターゲット予測)

- 機械学習

- N-gram featureの利用 (boostingなど)
- 木構造の索引 (tree kernel)

- 統計的機械翻訳

- フレーズアライメント
- 言語モデル

基礎知識

Suffix Arrays
Burrows Wheeler Transform

Suffix Arrays (SA) [Manber 1989]

● 入力: $T = t_1 t_2 t_3 \dots t_N$

● T の接尾辞(suffix): $S_k = t_k t_{k+1} t_{k+2} \dots t_N$

S_1	abraca\$		S_7	\$	7
S_2	braca\$		S_6	a\$	6
S_3	raca\$		S_1	abraca\$	1
S_4	aca\$	→	S_4	aca\$	→ 4
S_5	ca\$		S_2	braca\$	2
S_6	a\$		S_5	ca\$	5
S_7	\$		S_3	raca\$	3

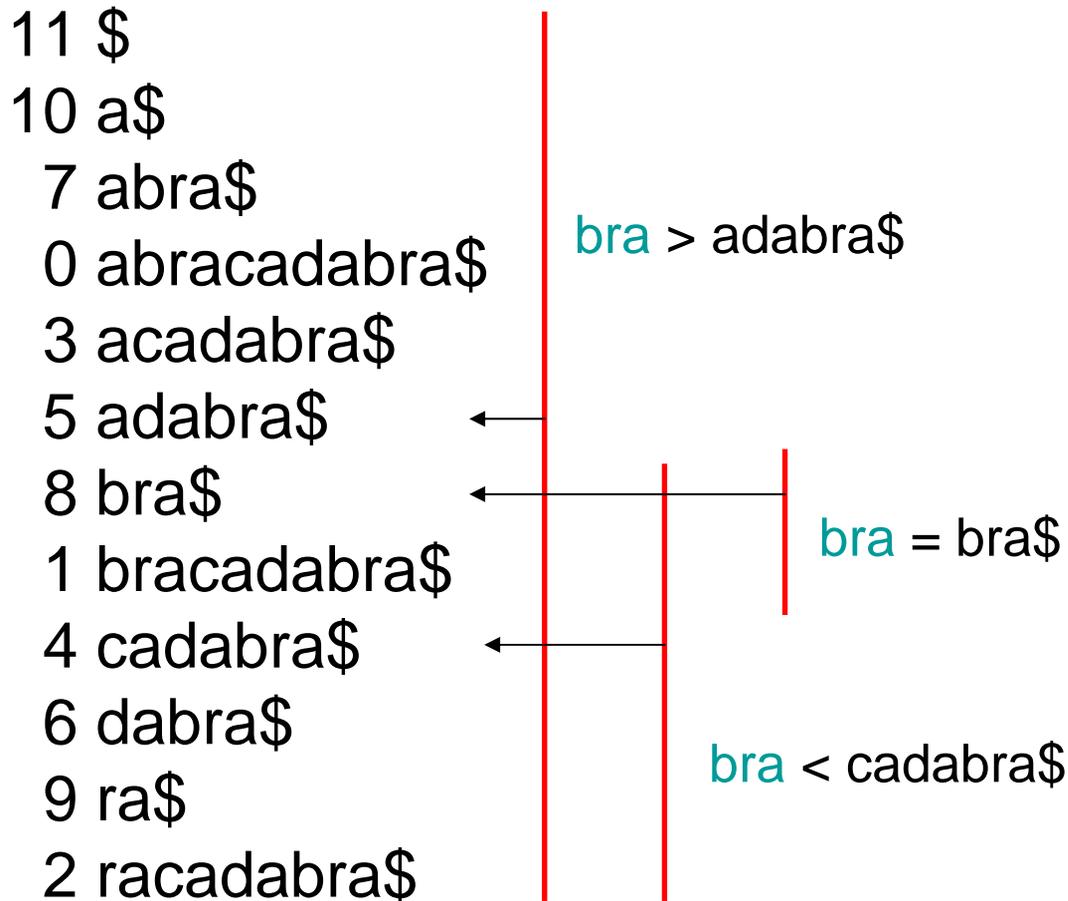
(1) T の全ての接尾辞を列挙

(3) 接尾辞の番号を抽出

(2) 接尾辞集合を辞書式順序でソートする

SAを使った検索

入力 $T = \text{abracadabra}\$$ パターン $P = \text{bra}$



時間計算量

$\text{occ}(P):$

$O(m \log n)$

$\text{loc}(P):$

$O(m \log n + \text{occ}(P))$

空間計算量

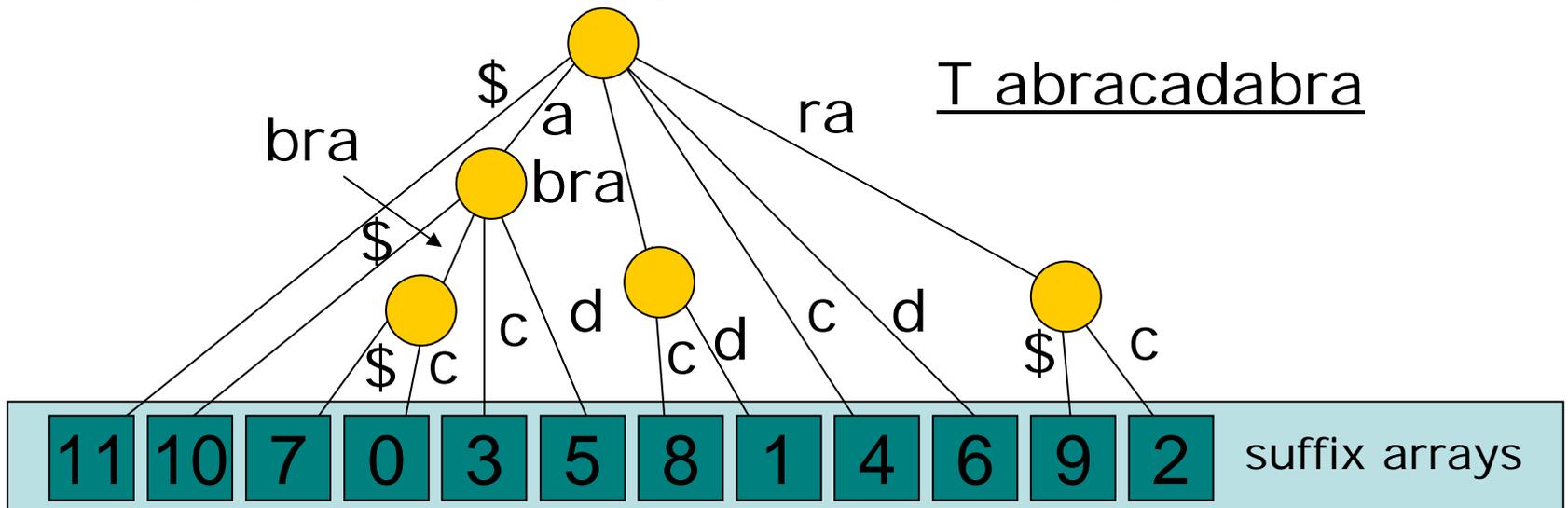
$\log n$ bit (4n byte)

Hgt配列を使うと

$\text{occ}(P)$ は $O(m + \log n)$

C.f. Suffix Trees (ST) [Weiner 1973]

- Tの全SuffixからなるTrieの圧縮(縮退)表現
- 多くの文字列アルゴリズムがST上で動く
- 非常に作業領域量が大い
- Compressed STは約 $9n$ bit
[Sadakane 2004] [Okanohara 2005]



Burrows Wheeler's Transform [1994] (BWT)

- 文字列に対する可逆変換
 - 最初はbzip2等の可逆圧縮に使われていた
- $BWT[i] := T[SA[i]-1]$
- `abracadabra$` \xrightarrow{BWT} `ard$rcaaaabb`
- BWTを適用したテキストは非常に圧縮しやすい
 - 同じ文脈の直前には同じ文字が現れやすい
 - t hese are great ...
 - t hese are possible ...
 - t hese were not of ..
 - t hese ...

例 BWT前

When Farmer Oak smiled, the corners of his mouth spread till they were within an unimportant distance of his ears, his eyes were reduced to chinks, and diverging wrinkles appeared round them, extending upon his countenance like the rays in a rudimentary sketch of the rising sun. His Christian name was Gabriel, and on working days he was a young man of sound judgment, easy motions, proper dress, and general good character. On Sundays he was a man of misty views, rather given to postponing, and hampered by his best clothes and umbrella : upon the whole, one who felt himself to occupy morally that vast middle space of Laodicean neutrality which lay between the Communion people of the parish and the drunken section, -- that is, he went to church, but yawned privately by the time the congregation reached the Nicene creed, - and thought of what there would be for dinner when he meant to be listening to the sermon. Or, to state his character as it stood in the scale of public opinion, when his friends and critics were in tantrums, he was considered rather a bad man ; when they were pleased, he was rather a good man ; when they were neither, he was a man whose moral colour was a kind of pepper-and-salt mixture. Since he lived six times as many working-days as Sundays, Oak's appearance in his old clothes was most peculiarly his own -- the mental picture formed by his neighbours in imagining him being always dressed in that way. He wore a low-crowned felt hat, spread out at the base by tight jamming upon the head for security in high winds, and a coat like Dr. Johnson's ; his lower extremities being

.....

BWTの重要な性質

- 同じ文字のBWT中に出現する順番はFでも同じ順番
- BWT, Fの順序はそれ
に続く文字列の順位で
決まる。どちらも同じ
文字列が順位決定に
使われる
- BW逆変換および、検
索でこの性質を利用

a ₁	\$	
r	a ₁	\$
d	a ₂	bra\$
\$	a ₃	bracadabra\$
r	a ₄	cadabra\$
c	a ₅	dabra\$
a ₂		bra\$
a ₃		bracadabra\$
a ₄		cadabra\$
a ₅		dabra\$
b		ra\$
b		racadabra\$

Legend: BWT, F

index	SA	BWT	F	Suffix
0	11	a ₁	\$	\$
1	10	r ₁	a ₁	a\$
2	7	d	a ₂	abra\$
3	0	\$	a ₃	abracadabra\$
4	3	r ₂	a ₄	acadabra\$
5	5	c	a ₅	adabra\$
6	8	a ₂	b ₁	bra\$
7	1	a ₃	b ₂	bracadabra\$
8	4	a ₄	c	cadabra\$
9	6	a ₅	d	dabra\$
10	9	b ₁	r ₁	ra\$
11	2	b ₂	r ₂	racadabra\$

BWT後のテキストから元のテキストを復元する。

BWTの性質(同じ文字の順番はBWT中でもFでも変わらない)を利用し、データをたどっていく。

F[i]に対応するBWT中の位置をLFmapping[i]とする。
例 LFmapping[3]=7
LFmapping[7]=11

\$の位置を求める i' = 3
a₃ を出力

a₃ のBWT中の位置を求め(7)
それをi'に代入 i' = 7

b₂ を出力

b₂ のBWT中の位置を求め(11)
それをi'に代入 i' = 11

..

逆BW变换 (LF-mapping)

```
void revBWT(char* bwt, int n){
    int count [0x100];
    memset(count,0,sizeof(int)*0x100);
    for (i = 0; i < n; i++) count[bwt[i]]++;
    for (int i = 1; i < 0x100; i++) count[i]+=count[i-1];
    int* LFmapping = new int[n];
    for (int i = n-1; i >= 0; i--){
        LFmapping[--count[bwt[i]]] = i;
    }
    int next = find(BWT,'$'); //find the position of '$'
    for (int i = 0; i < n; i++){
        next = LFmapping[next];
        putchar(bwt[next]);
    }
    delete[] LFmapping;
}
```

最近のSAの話題

- Compressed Suffix Arrays (CSA) 後で詳しく
- 高速な構築
 - 線形時間構築 [Ko 03] [Kim 03] [KS 03]
 - 現実的に一番速いものはこれらでは無い
c.f. dssort, msufsort, divsort
- メモリ効率の良い構築方法
 - CSAを直接構築 [sadakane 03]
 - CSAの高速構築 [mori 05] [okanohara 05]
 - 後で述べるwavelet treeを利用
- 近似マッチング [Huynh 05]

圧縮全文索引

Compressed Suffix Arrays (CSA)

[Grossi, Vitter 00][Sadakane 03][Grossi, Gupta, Vitter 03]

- 元のSAの作業領域量は $n \log n$ bit
 - SAは1からNまでの並び替えなので冗長性が無く圧縮はそのままでは不可能
- SAの代わりに次のように定義される SA^{-1} を保存
 - $[i] = SA^{-1}[SA[i] + 1]$
 - $SA^{-1}[i] = j$ s.t. $SA[j] = i$
 - SAをサンプリングした $SA_k[i] = SA[i \cdot k]$ も一緒に保存
- SA^{-1} を使ってSA上を移動。 SA_k からSAを求める
 - $SA[i] = SA[SA^{-1}[i]] - 1 = SA[SA^{-2}[i]] - 2 = \dots = SA[SA^{-n}[i]] - n$
- If $SA[SA^{-n}[i]] = p$ then $SA[i] = p - n$

I	SA	SA ⁻¹	$:= SA^{-1}[SA[i]+1]$	Suffix
0	11	3	3	\$
1	10	7	0	a\$
2	7	11	6	abra\$
3	0	4	7	abracadabra\$
4	3	8	8	acadabra\$
5	5	5	9	adabra\$
6	8	9	10	bra\$
7	1	2	11	bracadabra\$
8	4	6	5	cadabra\$
9	6	10	2	dabra\$
10	9	1	1	ra\$
11	2	0	4	racadabra\$

I	sampled SA	$:= SA^{-1}[SA[i]+1]$	Suffix
0		3	\$
1		0	a\$
2		6	abra\$
3	0	7	abracadabra\$
4		8	acadabra\$
5		9	
6	8	10	
7		11	
8	4	5	
9		2	
10		1	
11		4	racadabra\$

SA[7]

$$= SA[[7]] - 1 = SA[11] - 1$$

$$= SA[[11]] - 2 = SA[4] - 2$$

$$= SA[[4]] - 3 = SA[8] - 3$$

$$= 4 - 3$$

$$= 1$$

$[i] = SA^{-1}[SA[i] + 1]$ は一体何者？

- 直感的意味

- i 番目の suffix より 1 文字短い suffix の順序
- SA 表現では辞書隣接情報があり、位置隣接情報がない
位置表現では位置隣接情報があり、辞書隣接情報がない
- はそれら二つを結び付ける

- 定理:

もし $T[SA[i]] = T[SA[i+1]]$ ならば $[i] < [i+1]$

- 証明:

$T[SA[i]] = T[SA[i+1]]$ ならば $SA[i]$ と $SA[i+1]$ の順位はそれらより一文字短い suffix の順位で決定される (辞書式順序)

$$S_{SA[i]+1} < S_{SA[i+1]} \xrightarrow{\quad} SA^{-1}[SA[i]+1] < SA^{-1}[SA[i+1]+1]$$
$$\xrightarrow{\quad} [i] < [i+1]$$

I	SA	SA ⁻¹	$:= SA^{-1}[SA[i]+1]$	Suffix
0	11	3	3	\$
1	10	7	0	a\$
2	7	11	6	abra\$
3	0	4	7	abracadabra\$
4	3	8	8	acadabra\$
5	5	5	9	adabra\$
6	8	9	10	bra\$
7	1	2	11	bracadabra\$
8	4	6	5	cadabra\$
9	6	10	2	dabra\$
10	9	1	1	ra\$
11	2	0	4	racadabra\$

↓ 單調增加列

の圧縮

- は部分単調増加列なので圧縮可能
 - 差分列 d を次のように定義 $d[i] = a[i+1] - a[i]$
- delta符号を利用して $d[]$ を圧縮。サイズは $O(NH_0)$ [Sadakane 2003].
- Wavelet Treeを利用して d を圧縮。サイズは $O(NH_k)$ [Grossi 2003]
- 実用的には、 $a_k[i] = a[ik]$ と d を保存
 $a[i]$ を求めるには、 $t = i/k$ $r = i\%k$ であるとき
 $a[i] = a_k[tk] + d[tk] + d[tk+1] + \dots + d[tk+r-1]$

Self indexing property of CSA

[Sadakane 2003]

- CSA は検索時に元テキストを**必要としない!**
 - $T[i..j]$ を復元する $\text{substr}(i,j)$ もサポート
- Suffixの先頭文字は容易にわかる
 - 各文字の出現回数を保存しておけばよい
 - $D[i] :=$ テキスト中に実際に出現した文字
 - $C[i] :=$ $D[i]$ より小さい文字のテキスト中の出現回数
- $T[\text{SA}[i]] = D[j]$ s.t. $C[j] \leq i < C[j+1]$

```
void substr(int i, int j) {
    for (int p = SA-1[i]; i < j; i++, p = C[p]){
        int t = binarySearch(C);
        output(D[t]);
    }
}
```

I	SA	SA ⁻¹		Head of Suffix
0	11	3	3	\$
1	10	7	0	a
2	7	11	6	a
3	0	4	7	a
4	3	8	8	a
5	5	5	9	a
6	8	9	10	b
7	1	2	11	b
8	4	6	5	c
9	6	10	2	d
10	9	1	1	r
11	2	0	4	r

C[0] = 0

C[1] = 1

C[2] = 6

....

D[0] = '\$'

D[1] = 'a'

D[2] = 'b'

•substr(3,6)

p := SA⁻¹[3] = 4

decode[p]; // output 'a'

p = [p] = 8

decode[p]; // output 'c'

p = [p] = 5

decode[p]; // output 'a'

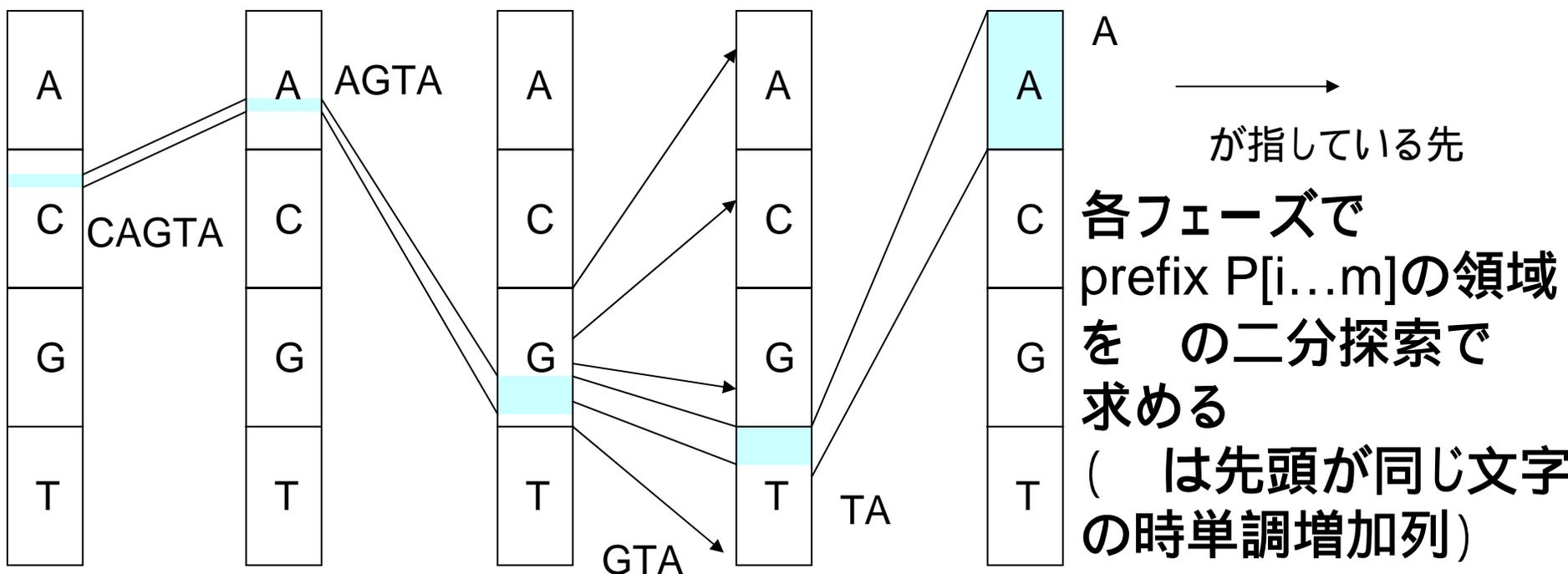
p = [p] = 9

decode[p]; // output 'd'

Backward Search

[Sadakane 2002] [Makinen 2004]

- 文字列探索時にSAを使うが、SAのlookupは計算量が多い。 だけを用いて探索が可能
- Search $P=CAGTA$ in backward ($P[m] P[m-1] \dots$)



FM-index [Ferragina 2000]

- もう一つの圧縮索引
 - 提案時は、実装が難しいとされていたが、最近現実的な実装方法が次々と見つかった
 - ちなみにLZ-indexという圧縮索引もある
- CSAと同様にocc, loc, substrをサポート
- BWTを基にしている
 - LF-mappingを検索に用いる

基本操作の定義

- rank (B, p, x)
 - B[0...p]中のxの出現回数を返す
- select (B, r, x)
 - B中でr回目のxの出現位置を返す
- 例
 - B=d#rcaaaabb
rank(B,6,a)=2
select(B,4,a)=7
select(B,2,b)=9
- rank, select操作を使ってoccを行う

$\text{occ}(P[1 \dots m], \text{BWT}[1 \dots n])$

$C[c]$: c より小さい文字の出現回数

BWT T のBWT後のテキスト

```
1.  i := m
2.  sp := 1; ep := n;
3.  while (sp < ep) and (i >= 1) do
4.      c := P[i];
5.      sp := C[c] + rank(BWT, c, sp-1)+1;
6.      ep := C[c] + rank(BWT, c, ep);
7.      i--;
8.  if (ep < sp)    return 0;
   else           return ep-sp+1;
```

I	SA	BWT	Head of Suffix
0	11	a ₁	\$
1	10	r ₁	a ₁
2	7	d	a ₂
3	0	\$	a ₃
4	3	r ₂	a ₄
5	5	c	a ₅
6	8	a ₂	b ₁
7	1	a ₃	b ₂
8	4	a ₄	c
9	6	a ₅	d
10	9	b ₁	r ₁
11	2	b ₂	r ₂

occ("ab", "ard\$rcaaaabb")

(1) "b"の領域を求める

I	SA	BWT	Head of Suffix
0	11	a ₁	\$
1	10	r ₁	a ₁
2	7	d	a ₂
3	0	\$	a ₃
4	3	r ₂	a ₄
5	5	c	a ₅
6	8	a ₂ ←	b ₁
7	1	a ₃ ←	b ₂
8	4	a ₄	c
9	6	a ₅	d
10	9	b ₁	r ₁
11	2	b ₂	r ₂

occ("ab", "ard\$rcaaaaabb")

(1) 'b'の領域を求める

(2) 現在の領域の直前である
(a₂ a₃)の領域を求める
(rank(BWT, 'a', ep),
rank(BWT, 'a'.sp))を求める

I	SA	BWT	Head of Suffix
0	11	a ₁	\$
1	10	r ₁	a ₁
2	7	d	a ₂
3	0	\$	a ₃
4	3	r ₂	a ₄
5	5	c	a ₅
6	8	a ₂	b ₁
7	1	a ₃	b ₂
8	4	a ₄	c
9	6	a ₅	d
10	9	b ₁	r ₁
11	2	b ₂	r ₂

FMcount("bra", "ard\$rcaaaaabb")

(1) 'b'の領域を求める

(2) 現在の領域の直前である
(a₂ a₃)の領域を求める
(rank(BWT, 'a', ep),
rank(BWT, 'a'.sp))

(3) "ab"の出現回数は
(3 - 2 + 1) = 2 回

CSA と FM-indexの関係

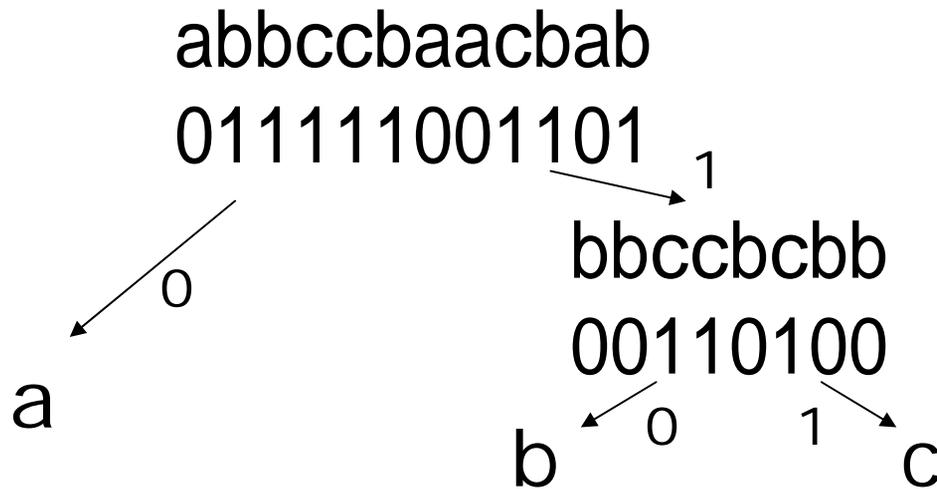
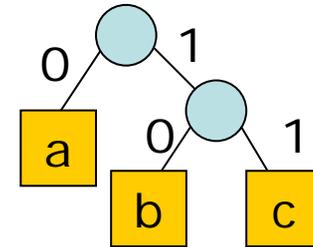
- CSAは SA^{-1} を、FM-indexは rank^{-1} を利用
 - $\text{rank}^{-1}[i] = \text{SA}^{-1}[\text{SA}[i]-1] = \text{C}[\text{T}[i]] + \text{rank}(\text{BWT}, \text{T}[i], i)$;
 - 別に逆を使っても問題ない
- 本質的に違うのはFM-indexが rank^{-1} を明示的に持たずにBWTのrankで実現しているところ
 - $\text{rank}(\text{BWT}, c, i)$ の効率良い実装は存在するか？
 - c が2値の場合は簡単。そうでない場合は困難
一般の c に対する実装方法が近年提案される

Rank, Select

- Bが2種類のアルファベットからなる場合
 - $\text{rank}(B,i,x)$, $\text{select}(B,i,x)$ を約 $1.5N$ bitの作業領域で $O(1)$ 時間で実現可能
 - Rankはテーブル参照 + popCount
 - Selectはrankの問題に帰着させる
- Bが3種類以上のアルファベットからなる場合
 - Wavelet Treeが現実的 [Grossi 2003]
 - rank, select を $O(NH_0)$ 作業領域、 $O(H_0)$ 時間で実現
 - 理論的には $O(NH_0)$ 作業領域、 $O(1)$ 時間で実現可能 [Ferragina 2005]

Example of Wavelet Tree

- $\Sigma = \{a, b, c\}$
 $a = 0_2$ $b = 10_2$ $c = 11_2$
- $T = \text{abbccbaacbab}$



最近の話題

Succinct*なデータ構造 Bit Array / Tree

* Succinct : n 個の要素を $O(n)$ (もしくは情報理論的な限界に漸近的に近い) 作業領域で各操作を定数時間で実行

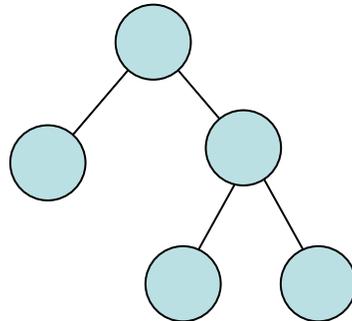
Succinct Bit Array

- Bit Vector $B[1 \dots N]$
 - rank (B, p, x) ($x=1,0$)
 - select (B, r, x) ($x = 1,0$)
 - 圧縮索引の様々な場面で利用
 - サンプルした位置が1、それ以外0のbit vector
- 定数時間rank/selectで作業領域は1.5Nbit程度
[Kim 2005]
- 定数時間select, $O(\log N)$ 時間rank, 作業領域 $O(H_0)$
[Okanoohara 2005]
 - 定数時間select/rankで作業領域 $O(H_0)$ は現在作業中

Balanced Parenthesis Tree (PT)

[Jacobson 85] [Munro 01] [Geary 05]

- 節点数+葉数= n の木構造を $2n$ bitの作業領域で表現
 - parent, first-child, siblingの操作時間は $O(1)$
 - c.f 情報理論的限界は $2n - o(n)$ bit
一般的に木はポインタを利用し $n \log n$ bitで表現
 - LISPと同様に木をDFSで辿った時、最初に辿った時('、出る時')'を出力した括弧列からなる

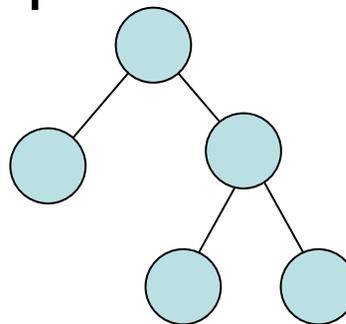


((() ((()))

0010010111

Balanced Parenthesis Tree (続)

- 括弧列上の次の操作を定数時間で行う
 - $\text{findopen}(x)$, $\text{findclose}(x)$: x に対応する括弧の位置を返す
 - $\text{enclose}(x)$: x を最もきつく囲む括弧対の開き括弧の位置を返す
- 木の操作は次のように実現
 - $\text{parent}(x) = \text{enclose}(x)$
 - $\text{sibling}(x) = \text{findclose}(x)+1$
 - $\text{first-child}(x) = x+1$



$(()()())$

0010010111

findcloseの例

- findclose(i)
- iの種類による場合分け
 - near
 - 前計算して作ったテーブル参照
 - pioneer
 - pioneerのみから作ったPTの答えを参照
 - far
 1. 直前のpioneer, pを探す (rankを利用)
 2. pまでの開括弧の数 - 閉括弧の数をtとする (rankを使う)
 3. pの対応する括弧対の位置からiの対応する括弧のブロックを求める
 4. pの位置とtから対応する括弧の位置をテーブル参照で求める

(((())))	(())
---	---	---	---	---	---	---	---	---	---	---	---

今後の展望

より複雑なデータ構造に対する索引

- 木、グラフの(圧縮)索引
 - 完全二分木に対するSuffix Tree [Shibuya 1999]
 - xbw [Ferragina 2006 to appear]
 - ラベル付き木に対するBWT
 - 作業領域量はtree entropy。Sub-pathなどが高速に求められる
- 近似マッチングの索引
 - ゲノムアライメントではpattern hunterという穴空きシードを使ったものが高速に近似マッチングを探索可能
 - ゲノムアライメント、統計的機械翻訳などで需要有り

理論的な枠組み

- , BWTに代わる表現は存在するか？
- 他の可逆変換の存在は？
- 理論的な限界の向上
 - NH_k bit の作業領域を用いて occ, loc を $O(m)$ 時間で実現可能か？