

## TODAY: Succinct data structures II (of 2)

- compact & succinct suffix arrays/trees
- $O(T \lg \lg T)$ -bit suffix array
- $O(T)$ -bit suffix array
- succinct suffix array  $\rightarrow$  tree

### Compact suffix arrays/trees:

TODAY \* Grossi & Vitter: [STOC 2000; SICALP 2005] first  
 $(1 + \frac{1}{\epsilon}) \underbrace{|T| \lg |\Sigma|}_{OPT^*} + \underbrace{2|T|}_{\text{if } |\Sigma| > 2} + O(\frac{T}{\lg \lg T})$  bits \*without compression

$O(\frac{P}{\log_{\Sigma} T} + \text{output}) \cdot \log_{\Sigma}^{\epsilon} T$  query

- Ferragina & Manzini: [FOCS 2000; JACM 2005]  
 $5 \cdot \underbrace{H_k(T)}_{k\text{th order empirical entropy}} \cdot |T| + O(T \frac{\epsilon + \lg T}{\lg \lg T} + T^{\epsilon} \Sigma^{\epsilon+1})$  Burrows-Wheeler Trans

for any fixed  $k$ <sup>↑</sup>  
 $= \sum_{|w|=k} \underbrace{\Pr\{w \text{ occurring}\}}_{= \# \text{ occurrences} / |T|} \cdot H_0(\text{string of successor characters of } w)$

$$= \sum_x P_x \lg \frac{1}{P_x} \cdot P_x = \frac{\# \text{ occ.}(wx)}{\# \text{ occ.}(w)}$$

= # bits/char. in OPT code depending on last  $k$  chars

$O(|P| + \text{output}) \cdot \lg^{\epsilon} T$  query

- Sadakane: [J. Alg. 2003] large alphabets

$\frac{1+\epsilon'}{\epsilon} \cdot H_0(T) \cdot |T| + O(T \lg \lg \Sigma + \Sigma \lg \Sigma)$  bits

$O(P \lg T + \text{output}) \cdot \lg^{\epsilon} T$  query for  $0 < \epsilon, \epsilon' < 1$

## Succinct:

- Grossi, Gupta, Vitter: [SODA 2003]  
 $H_k(T) \cdot |T| + O(T \lg \Sigma \cdot \frac{\lg \lg T}{\lg T})$  bits  
 $O(P \lg \Sigma + \frac{\lg^2 T}{\lg T} \lg \Sigma)$  query
- Ferragina, Manzini, Mäkinen, Navarro: [TALG 2007]  
 $H_k(T) \cdot |T| + O(T / \lg^\epsilon n)$  bits  
 $O(P + |\text{output}| \cdot \lg^{1+\epsilon} T)$  query +  $O(P)$  counting query

## Extras:

- low-space construction [Hon, Sadakane, Sung - Focs 2003;  
 $O(|T| \lg |\Sigma|)$  working Hon, Lam, Sadakane, Sung, Yu - Alg. 2007]
- suffix-tree ops. like maximal common substrings  
[Hon & Sadakane - CPM 2002; Sadakane - TCS 2007]
- document retrieval [Sadakane - J. Disc. Alg. 2007]
- dynamic [Chan, Hon, Lam, Sadakane - T. Alg. 2007]
- implementations & experiments [...]

# Compressed suffix arrays: [Grossi & Vitter - STOC 2000]

follow divide & conquer of linear-time suffix tree/array alg. [L16], but 2-way instead of 3-way

- start: text  $T_0 = T$

size  $n_0 = n$

suffix array  $SA_0 = SA$  in sorted order of suffixes

$SA[i] =$  index in  $T$  of  $i$ th suffix of  $T$

- step:  $T_{k+1} = \langle (T_k[2i], T_k[2i+1]) \text{ for } i=0, 1, \dots, n/2 \rangle$

$n_{k+1} = n_k/2 = n/2^k$

$SA_{k+1} = \frac{1}{2} \cdot$  extract even entries of  $SA_k$

- represent  $SA_k$  using  $SA_{k+1}$ :

① even-succ<sub>k</sub>( $i$ ) =  $\begin{cases} i & \text{if } SA_k[i] \text{ is even} \\ j & \text{if } SA_k[i] = SA_k[j] - 1 \text{ is odd} \end{cases}$

( $j$ th suffix starts right after  $i$ th suffix)

② even-rank<sub>k</sub>( $i$ ) = # even suffixes preceding  $i$ th suffix  
= # even values in  $SA_k[:i]$

③  $SA_k[i] = 2 \cdot SA_{k+1}[\text{even-rank}_k(\text{even-succ}_k(i))] - (1 - \text{is-even}_k(i))$

round to even suffix  $\Rightarrow$  in  $SA_{k+1}$

name of even suffix in  $SA_{k+1}$

index of even suffix in  $T_{k+1}$

index of even suffix in  $T_k$

unround

- stop recursion at level  $l = \lg \lg n \Rightarrow n_l = n / \lg n$

$\Rightarrow$  can afford naive ptr. encoding:  $n_l \lg n_l \leq n$  bits

$\Rightarrow O(\lg \lg n)$  query to  $SA[i]$ , given  $O(1)$ -time

is-even-suffix & even-succ & even-rank

$\text{is-even-suffix}_k(i) = \begin{cases} 1 & \text{if } SA_k[i] \text{ is even} \\ 0 & \text{else} \end{cases} \rightarrow n_k \text{ bits}$   
 $\text{even-rank}_k = \text{rank}_1 \text{ in is-even-suffix} \rightarrow O\left(n_k \frac{\lg \lg n_k}{\lg n_k}\right) \text{ bits}$  [L17]

even-succ<sub>k</sub>:

- trivial for  $SA_k[i]$  even ( $n_k/2$  such i's)
- store answers for  $SA_k[i]$  odd ( $n_k/2$  such i's)
- order by  $i \Rightarrow \text{even-succ}_k(i) = \text{odd-answers}[\text{odd-rank}_k(i)]$
- = order by odd suffix  $T_k[SA_k[i]:]$  =  $i - \text{even-rank}_k(i)$
- = order by  $(T_k[SA_k[i]], T_k[SA_k[i]+1:])$   
even
- = order by  $(T_k[SA_k[i]], T_k[SA_k[\text{even-succ}_k(i)]:])$  (def. of even-succ)
- = order by  $(T_k[SA_k[i]], \text{even-succ}_k(i))$  (SA<sub>k</sub> is a suffix array)  
answer!

- actually store these pairs, in order by value

- assume  $|\Sigma| = 2$  here

$\Rightarrow$  storing sorted array of  $n_k/2$  values,  $2^k + \lg n_k$  bits each

- store leading  $\lg n_k$  bits of each value  $v_i$  via unary differential encoding:  $0^{\text{lead}(v_1)} 1 0^{\text{lead}(v_2) - \text{lead}(v_1)} 1 \dots$

$\Rightarrow n_k/2$  1's &  $\leq 2^{\lg n_k} = n_k$  0's  $\Rightarrow \leq \frac{3}{2} n_k$  bits [L17]

-  $\text{lead}(v_i) = \text{rank}_0(\text{select}_1(i)) - \text{rank}_0(\text{select}_1(i-1))$

- store trailing  $2^k$  bits of each  $v_i$  explicitly in array

$\Rightarrow 2^k \cdot n_k/2 = n/2$  bits

$\Rightarrow \frac{1}{2} n + \frac{3}{2} n_k + O\left(\frac{n_k}{\lg \lg n_k}\right)$  bits

Total:  $\sum_{k=0}^{\lg \lg n} \left( n_k + \frac{1}{2} n + \frac{3}{2} n_k + O\left(\frac{n_k}{\lg \lg n_k}\right) \right)$   
rank & select

=  $\frac{1}{2} n \lg \lg n + 5n + O\left(\frac{n}{\lg \lg n}\right)$  bits - PROGRESS ~ BUT TOO BIG

# Compact suffix array: $O(n)$ bits

- store only  $1/\epsilon + 1$  levels of recursion:  
 $O_1 \epsilon l_1 2\epsilon l_2 \dots, l = \lg \lg n$
- represent  $SA_{k\epsilon l}$  using  $SA_{(k+1)\epsilon l}$  similarly:
  - "even"  $\rightarrow$  "divisible by  $2^{\epsilon l}$ " = "in  $SA_{(k+1)\epsilon l}$ "
  - is-even $_k = n_{k\epsilon l}$ -bit vector as before + rank $_1$  struct.
  - succ $_k(i) = j$  where  $SA_{k\epsilon l}[i] = SA_{k\epsilon l}[j] - 1$   
stored in  $n + 2n_{k\epsilon l} + O(\frac{n_{k\epsilon l}}{\lg \lg n_{k\epsilon l}})$  bits like even-succ  
 $\sim$  except can't skip  $1/2$  the values

- to compute  $SA_{k\epsilon l}[i]$ :

① follow succ $_k$  repeatedly until  $j$  in  $SA_{(k+1)\epsilon l}$

② recurse:  $SA_{(k+1)\epsilon l}[\text{rank}_1(j)]$   $\rightarrow$  is-even $_k$  bit vector

③ multiply by  $2^{\epsilon l}$ , subtract by # steps in ①

-  $\leq 2^{\epsilon l} = \lg^\epsilon n$  succ $_k$  calls per level

$\Rightarrow O(\lg^\epsilon n \lg \lg n) = O(\lg^\epsilon n)$  query to  $SA[i]$

- space:  $\sum_{k=0}^{1/\epsilon} (n_{k\epsilon l} + n + 2n_{k\epsilon l} + O(\frac{n_{k\epsilon l}}{\lg \lg n_{k\epsilon l}}))$

$$= (\frac{1}{\epsilon} + 6)n + O(\frac{n}{\lg \lg n}) \text{ bits}$$

- optimizations:

- succ $_k$  free at level  $k=0 \Rightarrow \sum \leq 4n_{\epsilon l} = O(\frac{n}{\lg^\epsilon n})$

- store is-even $_k$  as succinct dictionary (with rank)

$$\Rightarrow \lg \binom{n_{k\epsilon l}}{n_{(k+1)\epsilon l}} + o \sim n_{(k+1)\epsilon l} \lg \frac{n_{k\epsilon l}}{n_{(k+1)\epsilon l}} + o = O(n \frac{\lg \lg n}{\lg^\epsilon n})$$

$$\Rightarrow (\frac{1}{\epsilon} + 1) \cdot n + O(\frac{n}{\lg \lg n}) \text{ bits}$$

$\rightarrow$  [Brodnik & Munro] of L17

**OPEN:**  $O(n)$  bits &  $o(\lg^\epsilon n)$  query

## Compact suffix tree given suffix array:

[Munro, Raman, Rao - J. Alg. 2001]

- store compressed binary trie on  $2n+1$  nodes in  $4n + o(n)$  bits using balanced parens. [L17]

↳ "skipping the skips": can't store edge lengths

- during search for  $P$ , maintain letter depth of node: to descend  $(x) \rightarrow (y)$ :

- compute length of edge by finding longest match between leftmost-leaf( $y$ ) & rightmost-leaf( $y$ ),

starting at letter depth of  $x$  (leaf rank + SA)

- check that  $P$  matches these letters

$\Rightarrow O(|P| + |\text{output}|) \cdot \text{cost for SA lookup}$

↳ # matches via # leaves in subtree

↳ enumerate  $k$  matches via leaf select + rank

Grossi & Vitter achieve  $O\left(\frac{P}{\log_{\Sigma} T} + |\text{output}| \cdot \log_{\Sigma}^{\epsilon} T\right)$   
↳ word RAM

### Leaf ops:

leaf =  $()$

leaf-rank =  $\text{rank}()$  (here)

leaf-select =  $\text{select}()$  ( $i$ )

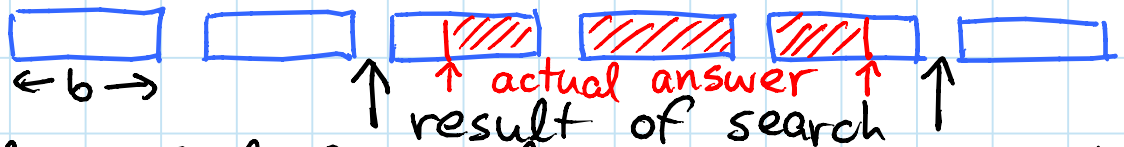
# leaves in subtree =  $\text{rank}()$  (matching) of parent  
-  $\text{rank}()$  (here)

leftmost-leaf =  $\text{leaf-select}(\text{leaf-rank}(\text{here}) + 1)$

rightmost-leaf =  $\text{leaf-select}(\text{leaf-rank}(\text{matching of parent}) - 1)$

# "Succinct" suffix tree, given suffix array [MRR'01]

- store answer to every query with  $|P| < b$   
 $\Rightarrow O(2^b \lg n)$  bits =  $O(\sqrt{n})$  if  $b \leq \frac{1}{3} \lg n$
- now consider  $|P| \geq b$
- use suffix tree above on every  $b$ th suffix in suffix order (keep every  $b$ th leaf in tree)
- search in tree narrows to interval of size- $b$  blocks in SA:



$\Rightarrow$  need to find first & last match in a block

- lookup table: (not really needed)

for any  $b$   $b$ -bit strings in sorted order  
& any  $\leq b$ -bit query string.

find first & last prefix match

$\Rightarrow O(2^{b^2+b} \lg b)$  bits =  $O(\sqrt{n})$  if  $b \leq \frac{1}{5} \sqrt{\lg n}$

- to search in a block:

① locate the  $b$  suffixes in  $T$  using SA

① read next  $b$  bits from  $P$

& from all  $b$  suffixes (in  $T$ )

& use lookup table to narrow range

② repeat  $\lceil |P|/b \rceil$  times

$\Rightarrow O(P)$  time +  $O(b)$  · cost of SA query

$\Rightarrow O(P + \lg^\epsilon n)$  time for  $b = O(\lg^\epsilon n)$

$\Rightarrow$  dominated by  $O(P \cdot \text{SA query})$  from previous struct.

-  $O(\frac{n}{b})$  bits plus size of suffix array

$2^{b^2}$   
 $2^b$   
 $O(\lg b)$

$O(b)$  ·  
SA query  
+  
 $O(b)$   
·  $O(P/b)$