

## Compact suffix arrays/trees:

<sup>TODAY</sup>\* Grossi & Vitter: [STOC 2000; SICOMP 2005]

$$\left(\frac{1}{\varepsilon} + O(1)\right) \underbrace{|T| \lg |\Sigma|}_{OPT} \text{ bits}$$

first

$$O\left(\frac{|P|}{\log_{|\Sigma|} |T|} + |\text{output}| \cdot \log_{|\Sigma|}^{\varepsilon} |T|\right) \text{ query}$$

- Ferragina & Manzini: [FOCS 2000] compressed  
 $5 \cdot H_k(T) \cdot |T| + O\left(\frac{|T|}{\log |T|}(|\Sigma| + \lg \lg |T|) + |T|^{\varepsilon} |\Sigma|^{O(|\Sigma|)}\right)$  bits  
 $\underbrace{H_k(T)}_{k\text{th order empirical entropy}}$  for any fixed  $k$   
 $= \sum_{|w|=k} \Pr\{w \text{ occurring}\} \cdot H_0(\text{string of successor characters of } w)$   
 $= \# \text{ bits/char. in OPT code depending on last } k \text{ chars}$

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

- Sadakane: [J. Alg. 2003] large alphabets  
 $\frac{1}{\varepsilon} \cdot H_0(T) \cdot |T| + O(|T| \lg \lg |\Sigma| + |\Sigma| \lg |\Sigma|)$  bits  
 $O(|P| \lg |T| + |\text{output}| \cdot \log^{\varepsilon} |T|) \text{ query}$

- Grossi, Gupta, Vitter: succinct [SODA 2003]  
 $H_k(T) \cdot |T| + O(|T| \lg |\Sigma| \cdot \frac{\lg \lg |T|}{\lg |T|})$  bits  
 $O(|P| \lg |\Sigma| + \lg^{O(1)} |T|) \text{ query}$

- 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]

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

follow divide & conquer of linear-time suffix tree/array alg. [L7], 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

$SA[i] = \text{index in } T \text{ of } i\text{th suffix of } T$

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

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

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

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

$$\textcircled{1} \quad \underline{\text{even-succ}_k(i)} = \begin{cases} i & \text{if } SA_k[i] \text{ is even} \\ j & \text{if } SA_k[i] = SA_k[j]-1 \text{ is odd} \end{cases}$$

$$\textcircled{2} \quad \underline{\text{even-rank}_k(i)} = \# \text{even suffixes preceding } i\text{th suffix} \\ = \# \text{even values in } SA_k[:i]$$

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

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  $l = \lg \lg n \Rightarrow n_l = \frac{n}{\lg n}$

$\Rightarrow$  can afford naive  $n_l \lg n_l \leq n$ -bit ptr. encoding

$\Rightarrow O(\lg \lg n)$  query to  $SA[i]$ , given  $O(1)$ -time  
is-even-suffix & even-succ & even-rank

is-even-suffix<sub>k</sub>(i) =  $\begin{cases} 1 & \text{if } SA_k[i] \text{ is even} \\ 0 & \text{else} \end{cases}$   $\rightarrow n_k$  bits  
even-rank<sub>k</sub> = rank<sub>1</sub> in is-even-suffix  $\rightarrow O(n_k \frac{\lg \lg n_k}{\lg n_k})$  bits

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  $\Rightarrow$  even-succ<sub>k</sub>(i) = odd-answers[odd-rank<sub>k</sub>(i)]
- = order by odd suffix  $T_k[SA_k[i]:]$   $= i - \text{even-rank}_k(i)$
- = order by  $(T_k[SA[i]], T_k[\underbrace{SA_k[i]+1:}_{\text{even}}])$
- = order by  $(T_k[SA[i]], T_k[SA_k[\text{even-succ}_k(i)]:])$
- = order by  $(T_k[SA[i]], \underbrace{\text{even-succ}_k(i)}_{\text{answer!}})$  ( $SA_k$  is a suffix array)

(def. of even-succ)

- 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:  $0^{\text{lead}(v_1)} 1 0^{\text{lead}(v_2) - \text{lead}(v_1)} 1 \dots$

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

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

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

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

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

Total:  $\sum_{k=0}^{\lg \lg n} \left( n_k + \frac{1}{2}n + \frac{3}{2}n_k + O\left(\frac{n_k}{\lg \lg n_k}\right) \right)$   
 $= \frac{1}{2}n \lg \lg n + 5n + O\left(\frac{n}{\lg \lg n}\right) \text{ bits}$

TOO BIG

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

- store only  $\frac{1}{\varepsilon} + 1$  levels of recursion:  
 $O_1 \varepsilon l, 2\varepsilon l, \dots, l = \lg \lg n$
- represent  $SA_{k\ell}$  using  $SA_{(k+1)\ell}$ . Similarly:
  - "even"  $\rightarrow$  "divisible by  $2^{\ell}$ " = "in  $SA_{(k+1)\ell}$ "
  - $is-even_k = n_{k\ell}$  -bit vector as before + rank<sub>1</sub> struct.
  - $succ_k(i) = j$  where  $SA_{k\ell}[i] = SA_{k\ell}[j] - 1$   
 stored in  $n + 2n_{k\ell} + O(\frac{n_{k\ell}}{\lg \lg n_{k\ell}})$  bits like even-succ
- to compute  $SA_{k\ell}[i]$ :
  - ① follow  $succ_k$  repeatedly until  $j$  in  $SA_{(k+1)\ell}$
  - ② recurse:  $SA_{(k+1)\ell}[\underline{\text{rank}_1(j)}]$   $\xrightarrow{\text{is-even}_k \text{ bit vector}}$
  - ③ multiply by  $2^{\ell}$ , subtract by # steps in ①
- $\leq 2^{\ell} = \lg^\varepsilon n$  succ<sub>k</sub> calls per level
- $\Rightarrow O(\lg^\varepsilon n \lg \lg n) = O(\lg^\varepsilon n)$  query to  $SA[i]$
- Space:  $\sum_{k=0}^{\frac{1}{\varepsilon}} (n_{k\ell} + n + 2n_{k\ell} + O(\frac{n_{k\ell}}{\lg \lg n_{k\ell}}))$   
 $= (\frac{1}{\varepsilon} + 6)n + O(\frac{n}{\lg \lg n})$  bits
- optimizations:

$- succ_k$  free at level  $k=0 \Rightarrow \sum \leq 4n_{\ell} = O(\frac{n}{\lg^\varepsilon n})$   
 $- store is-even_k as succinct dictionary (with rank)$   
 $\Rightarrow \lg(n_{k\ell}) + o \sim n_{(k+1)\ell} \lg \frac{n_{k\ell}}{n_{(k+1)\ell}} + o = O(n \frac{\lg \lg n}{2^{\ell}})$

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

OPEN:  $O(n)$  bits &  $O(\lg^\varepsilon n)$  query

## Compact suffix tree given suffix array:

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

- store compressed binary trie on  $2n+1$  nodes in  $4n \log(n)$  bits using balanced parens. [L22]
    - ↪ "skipping the skips": can't store edge lengths
  - during search for  $P$ , maintain letter depth of node: to descend  $\xrightarrow{x} \circlearrowright y$ :
    - 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
- $\Rightarrow O((|P| + \text{loutput}) \cdot \text{cost for SA lookup})$

↪ # matches via # leaves in subtree

↪ enumerate  $k$  matches via leaf select+rank

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

### Leaf ops:

leaf = ()

leaf-rank = rank<sub>(())</sub>(here)

leaf-select = select<sub>(())</sub>(i)

# leaves in subtree = rank<sub>(())</sub>(matching) of parent  
- rank<sub>(())</sub>(here)

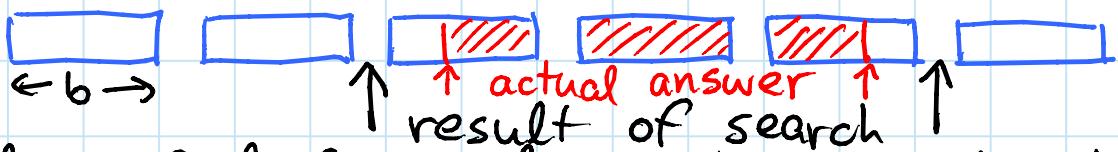
leftmost-leaf = leaf-select(leaf-rank(here) + 1)

rightmost-leaf = leaf-select(leaf-rank(matching) of parent - 1))

# "Succinct" suffix tree, given suffix array

[Munro, Raman, Rao: fixed]

- 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:



⇒ need to find first & last match in a block

- lookup table:

for any  $b$   $b$ -bit strings in sorted order

& any  $\leq b$ -bit query string,

find first & last prefix match

$$\begin{aligned} & \alpha^{2^b} \\ & 2^b \\ & O(\lg b) \end{aligned}$$

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

- to search in a block:

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

$$O(b)^+$$

① read next  $b$  bits from  $P$

$$(O(1))^+$$

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

$$(O(b))$$

② repeat  $\lceil \frac{|P|}{b} \rceil$  times

$$\cdot O(P/b)$$

$$\Rightarrow O(|P| + b) \text{ time} = O(|P| + \lg \lg n) \text{ if } b = O(\lg \lg n)$$

*forgotten in [MRR]*

× cost of SA query  $\Rightarrow O(|P| \lg^\epsilon n)$

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