

Leaf operations:

leaf

rightmost child
of parent

)()

leaf rank

rank₎₍₎(here)

leaf count
in subtree

rank₎₍₎[matching] of parent]
— rank₎₍₎(here)

leaf select

select₎₍₎(i)

leftmost leaf
in subtree

leaf select(1+

leaf rank(here))

rightmost leaf
in subtree

leaf select(leaf rank(
char. before matching)
of parent)))

Compact suffix arrays/trees

* Grossi & Vitter [STOC 2000; SICOMP 2005]:
 TODAY $(\frac{1}{\varepsilon} + O(1)) \underbrace{|T| \lg |\Sigma|}_{OPT}$ bits

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

- Ferragina & Manzini [FOCS 2000]:

$$\underbrace{5 \cdot H_k(T) \cdot |T|}_{k\text{-order "compressed OPT"}} + O\left(\frac{|T|}{\log|T|}(|\Sigma| + \lg \lg |T|) + |T|^{\varepsilon} |\Sigma| 2^{|\Sigma| \lg |\Sigma|}\right)$$

bits, for any fixed k

- $H_k(T)$ = k th order empirical entropy

$$= \sum_{\substack{w \in \Sigma \\ |w|=k}} \Pr\{w \text{ occurring}\} \cdot H_0(\text{string of successor characters of } w)$$

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

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

- Sadakane [J. Alg. 2003]

$$\frac{1}{\varepsilon} \cdot H_0(T) \cdot |T| + O(|T| \lg \lg |\Sigma| + |\Sigma| \lg |\Sigma|) \text{ bits}$$

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

- Grossi, Gupta, Vitter [SODA 2003]: succinct

$$H_k(T) \cdot |T| + O(|T| \lg |\Sigma| \cdot \frac{\lg \lg |T|}{\log|T|}) \text{ bits}$$

$$O(|P| \lg |\Sigma| + \lg^{O(1)}|T|) \text{ query}$$

- low-space construction: [Han, Sadakane, Sung - FOCS 2003]

$$O(|T| \lg |\Sigma|) \text{ bits working space}$$

$O(|T|)$ time for suffix array; $O(|T| \lg^{\varepsilon}|T|)$ for ^{suffix} tree

- more "suffix-tree like" [Han & Sadakane - CPM 2002]

e.g. find maximal common substrings

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

- follow divide & conquer of linear-time suffix tree alg.
but 2-way instead of 3-way
- start: text $T_0 = T$
size $n_0 = n$
suffix array $SA_0 = SA$:
 $SA[i] = \text{index in } T \text{ of } i\text{th suffix of } T$
in sorted order of suffixes of T
- step: text $T_{k+1} = \langle (T_k[2i], T_k(2i+1)) \text{ for } i=0, 1, \dots, \frac{n}{2} \rangle$
 size $n_{k+1} = \frac{n_k}{2} = \frac{n}{2^k}$
 suffix array SA_{k+1}
 $= \frac{1}{2} \cdot \text{extract even entries of } SA_k$
- represent SA_k using SA_{k+1} :
 - ① even-succ_k(i) = $\begin{cases} i & \text{if } SA_k[i] \text{ is even} \\ j & \text{if } SA_k[i] = SA_k[j] - 1 \text{ odd} \end{cases}$
 - ② even-rank_k(i) = # even suffixes preceding i th suffix
 $= \# \text{ even values in } SA_k[:i]$
 - ③ $SA_k[i] = 2 \cdot SA_{k+1}[\underbrace{\text{even-rank}_k(\text{even-succ}_k(i))}_{\text{name of even suffix in } SA_{k+1}}] - \underbrace{\{1 \text{ if } \frac{SA_k[i]}{\text{odd}}\}}_{\text{round to even suffix } \Rightarrow \text{in } SA_{k+1}}$
round to even suffix \Rightarrow in SA_{k+1}
index of even suffix in T_{k+1}
index of even suffix in T_k round back
- 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 encoding
 $\Rightarrow O(\lg \lg n)$ query to $SA[i] \dots$ if even-succ & even-rank & is-even-suffix are $O(1)$.

is-even-suffix_k(i) = $\begin{cases} 1 & \text{if } SA_k[i] \text{ is even} \\ 0 & \text{else} \end{cases}$ $\sim n_k$ -bit vector

even-rank_k = rank₁ in is-even-suffix $\sim O(n_k \frac{\lg \lg n_k}{\lg n_k})$ bits

even-succ_k:

- 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_k(i) = odd-answers[odd-rank_k(i)]
- = order by odd suffix $T_k[SA_k[i]:]$ $i - \text{even-rank}_k(i)$
- = order by $(T_k[SA_k[i]], T_k[SA_k[i]+1:])$
- = order by $(T_k[SA_k[i]], T_k[SA_k[\text{even-succ}_k(i)]:]^{\text{even}})$
- = order by $(T_k[SA_k[i]], \underbrace{\text{even-succ}_k(i)}_{\text{answer}})$ (SA_k is suffix array)
- actually store these pairs \sim in order by value assuming $|S|=2$
- \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 by 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
- random access: $\text{leading}(v_i) = \text{rank}_0(\text{select}_1(i))$
- store trailing 2^k bits of each v_i explicitly in array
- $\Rightarrow 2^k \cdot \frac{n_k}{2} = n/2$ bits
- $\Rightarrow \frac{1}{2}n + \frac{3}{2}n_k + O\left(\frac{n_k}{\lg \lg n_k}\right)$ bits

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

Compact suffix array: $O(n)$ bits [Grossi & Vitter]

- store only $\frac{1}{\varepsilon} + 1$ levels of recursion:

$$0, \varepsilon l, 2\varepsilon l, \dots, l = \lg \lg n$$

- represent $SA_{k\ell}$ using $SA_{(k+1)\ell}$ similarly:

- "even" now means "even εl times" (in $SA_{(k+1)\ell}$)

- is-even_k: $n_{k\ell}$ -bit vector as before + rank₁ structure

- 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 is in $SA_{(k+1)\ell}$

② recurse: $SA_{(k+1)\ell}[\underline{\text{rank}_1(j)}]$ is-even_k bit vect.

③ multiply by $2^{\varepsilon l}$, subtract by # calls to succ_k

\Rightarrow search cost $\leq 2^{\varepsilon l} = \lg^\varepsilon n$ succ_k calls per level

$$= O(\lg^\varepsilon n \lg \lg n) = O(\lg^\varepsilon n) \text{ time}$$

- space: $\sum_{k=0}^{\frac{1}{\varepsilon}} [n_{k\ell} + O(n_{k\ell} \frac{\lg \lg n_{k\ell}}{\lg n_{k\ell}}) + n + 2n_{k\ell} + O(\frac{n_{k\ell}}{\lg \lg n_{k\ell}})]$
 $= 7(\frac{1}{\varepsilon} + 1)n + O(\frac{n}{\lg \lg n})$ bits

- optimizations:

- succ_k free at level k=0 $\Rightarrow \leq 4n_{\ell} = O(\frac{n}{\lg^\varepsilon n})$ bits

- store is-even bit vector as succinct dictionary (with rank)

$$\Rightarrow \underbrace{\lg \binom{n_{k\ell}}{n_{(k+1)\ell}}}_{\approx n_{(k+1)\ell} \lg \frac{n_{k\ell}}{n_{(k+1)\ell}}} + O(n_{k\ell} \frac{(\lg \lg n_{k\ell})^2}{\lg n_{k\ell}}) \text{ bits}$$

$$\approx n_{(k+1)\ell} \lg \frac{n_{k\ell}}{2^{\varepsilon l}} = O(n \frac{\lg \lg n}{\lg^\varepsilon n}) \text{ bits}$$

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

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 $4n+o(n)$ bits using balanced parens.
"skipping the skips" — don't know edge lengths
 - as search for P proceeds, maintain letter depth of node
 - to descend $\xrightarrow{x} \downarrow y$:
 - 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
- $\Rightarrow O((|P| + \underline{\text{output}}) \cdot \text{cost for SA lookup})$
- ↳ #matches via leaf size; enumerate k via leaf ^{select} & rank

"Succinct" suffix tree, given suffix array [Munro, Raman, Rao]

- use structure above on every $b = \frac{1}{2}\sqrt{\lg n}$ suffixes in suffix order (i.e. keep every b th leaf of suffix tree)
- search narrows to an 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, for any $\leq b$ -bit query string, find first & last prefix match
 $\Rightarrow O(2^{b^2} 2^b \lg b) = O(\sqrt[4]{n} 2^{\frac{1}{2}\sqrt{\lg n}} \lg \lg n) = O(n^{\frac{1}{4}+\varepsilon})$ bits
- to search in a block:
 - ① read next b bits from all b suffixes in block & from P via T
 $\Rightarrow O(b)$ time
 - ② repeat $\lceil \frac{|P|}{b} \rceil$ times
 $\Rightarrow O(|P| + b)$ time at bottom
↳ reducible to $\lg^\varepsilon n$ by reducing b
- $O(\frac{n}{b})$ bits plus size of suffix array

Grossi & Vitter achieve $O(\frac{|P|}{\log_{|\Sigma|}|T|} + |\text{output}| \cdot \log_{|\Sigma|}^\varepsilon |T|)$ query & $O(|T|)$ bits