

Suffix arrays

DS 11/17

DS for string algorithms in text processing, pattern matching, data compression, bioinformatics, ... (any linear data, e.g. GPS record)

substring search problem

find all occurrences of a substring β in a given string α , $n=|\alpha|$, $m=|\beta|$

algs: Robin-Karp (rolling hashing as a filter) ... $\mathcal{O}(n \cdot m)$ worst case time

Knuth-Morris-Pratt } build DS for β , use for any d ... $\mathcal{O}(n+m)$ time
Aho-Corasick } (partial match table, finite automaton)

DS: build DS for a fixed α , use for any β .

suffix array, build (precomputation) $\mathcal{O}(n)$ time \Rightarrow query in $\mathcal{O}(\log n + m)$
(suffix tree) (we show only weaker results) (if only count reported)

Notation:

$\alpha = \alpha[0] \dots \alpha[n-1]$, $n=|\alpha|$ (symbols labeled from 0)

$\alpha[i:j] := \alpha[i] \dots \alpha[j-1]$ a substring of length $j-i$ if $i < j$,
 ϵ (empty string) if $i \geq j$.

$\alpha[:j]$ prefix of length j in α , $0 \leq j \leq n$ (i.e. $\alpha[0:j]$)

$\alpha[i:]$ suffix starting in $\alpha[i]$, $0 \leq i \leq n$ (i.e. $\alpha[i:n]$)

$\alpha \leq \beta$ lexicographical order: α is a prefix of β or there is i s.t.
 $\alpha[i] < \beta[i]$ and $\alpha[:i] = \beta[:i]$.

\uparrow fixed total order on alphabet

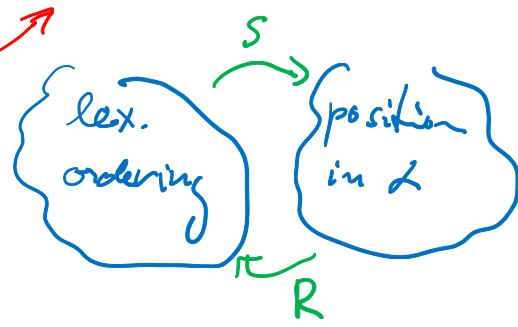
Note - To avoid this, we may use "padding" $\alpha\$, \beta\%$ where $\$$ (end of string) is a new symbol (smaller than other symbols).

• We have no assumption on the alphabet, is finite but can be arbitrarily large.

ex: $\alpha = \text{bananas}$, $n = |\alpha| = 7$

$\beta = \text{na}$, $m = |\beta| = 2$
occurrences of β in α

i	$\alpha[S[i]:]$	$S[i]$	$R[i]$	$L[i]$
0	E	7	4	0
1	ananas	1	1	3
2	anas	3	5	1
3	as	5	2	0
4	bananas	0	6	0
5	anas	2	3	2
6	nas	4	7	0
7	s	6	0	-



in lex order

all suffixes in lex. order

where the lex. suffix starts ind

longest common prefix with the next row
lex. rank (which row) of $\alpha[i:]$

Def: The **suffix array** for a string α is a permutation $S[0..n]$ s.t. $\alpha[S[i]:] < \alpha[S[i+1]:]$ for every $0 \leq i < n$, i.e. $S[i]$ is the starting position of the i -th lex. smallest suffix in α .

☺ All occurrences of β in α form an interval in the suffix array. (as prefixes of suffixes)

⇒ All occurrences of β in α can be enumerated from S by a binary search in time $O(n \log n + p)$, $p = \# \text{ occurrences}$.
(to be improved) → not needed for count only

Def: The **rank array** $R[0..n]$ is an inverse permutation to S , i.e. $R[i]$ is # lex. smaller suffixes than $\alpha[i:]$.

Def: The **LCP array** (longest common prefix) $L[0..n-1]$ is given by:
 $L[i] := \text{LCP}(\alpha[S[i]:], \alpha[S[i+1]:])$ for every $0 \leq i < n$,
where $\text{LCP}(p, \beta) := \max. k \text{ s.t. } p[1:k] = \beta[1:k]$.

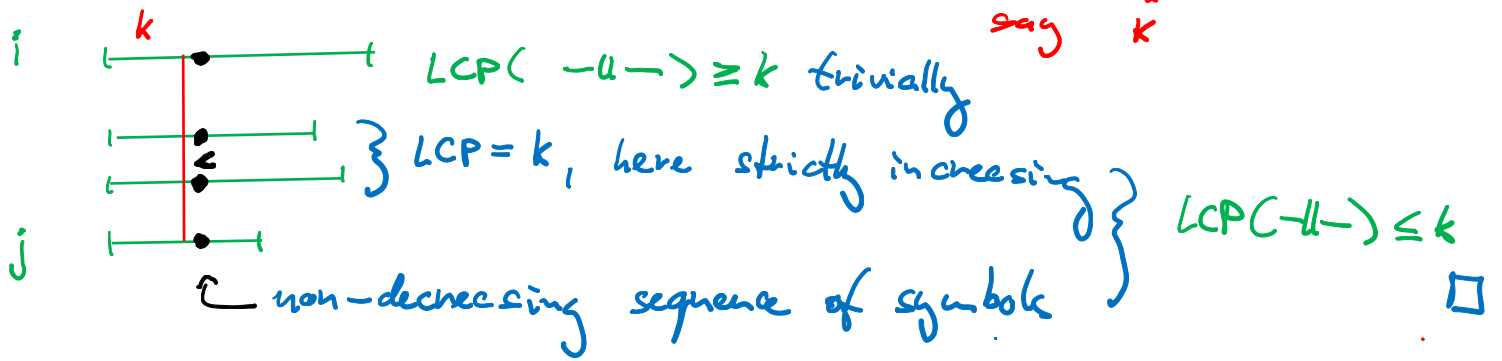
LCP of any two suffixes

☺ For any $0 \leq i < j \leq n$

minimum of "adjacent steps"

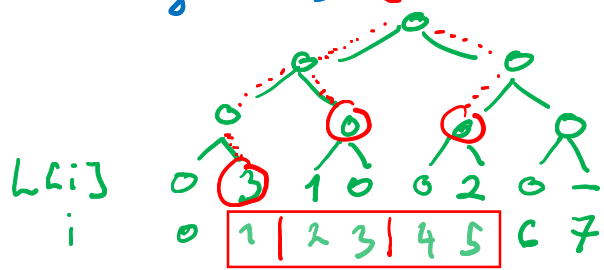
$LCP(\alpha[S[i]:], \alpha[S[j]:]) = \min(L[i], \dots, L[j-1])$

Pf:



RMQ (range minimum queries) for L (to find min in range [i:j])

ex: 1D-range trees (next lecture)



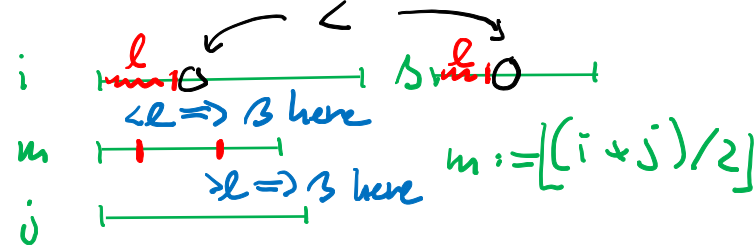
← minimum in inner nodes (subintervals)

$\Rightarrow O(\log n)$ query time
 $O(n)$ build time

Note: There is a DS with $O(1)$ query time, $O(n)$ build time. (without proof)

☺ Given LCP for any two suffixes, the lex. smallest suffix of α with prefix β can be found in $O(m + \log n)$ time. (or largest)

idea: current range $[i, j]$, assume $l = LCP(\beta, \alpha[S[i]:])$ known.



if $LCP(\alpha[S[l]:], \alpha[S[m]:]) < l$, continue in $[i, m]$

if $* > l$, continue in $[m, j]$

if $* = l$, compare β with $\alpha[S[m]:]$ from $(l+1)$ -th symbol.

\Rightarrow uses $\leq m$ comparisons in β

\Rightarrow All occurrences of β in α can be enumerated in $O(m + \log n + p)$ time, $p = \#$ occurrences.

Applications of LCP

substrings of length k

1) histogram of k-grams = # occurrences of each k-gram

ex:

i	$L[S[i:]]$	$L[i]$	$k=2$
0	ϵ	0	← too short
1	ananas	3] an
2	anas	1	
3	as	0	
4	bananas	0] ba
5	nanas	2] na
6	nas	0	
7	s	-	← too short

\Rightarrow ranges $[i, j]$ in LCP array s.t. $L[i-1] < k$, $L[i] \geq k \neq i \leq j < i+1$, $L[j] < k$

\Rightarrow single pass through L in $O(n)$ time

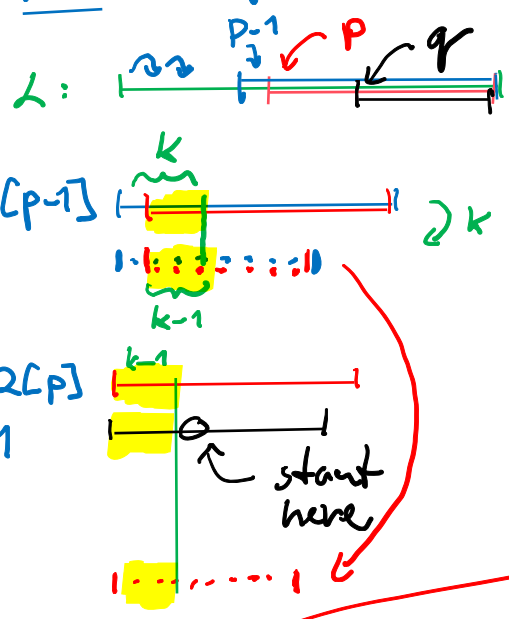
2) longest repeating substring in s = $\max_i L[i] \Rightarrow O(n)$ time
 ("redundance" in s)

3) longest common substring of s and t : from s from t
 LCP for $s \neq t$, then $\max L[i]$ s.t. $S[i] < |s|$ and $S[i+1] > |s|$ or vice versa
 ↑ separator

Construction of LCP array (Kasai's alg.)

given S, R construct L in $O(n)$ time

idea: if $k = L[R[p-1]]$ then $L[R[p]] \geq k-1$.



BUILD LCP

```

k ← 0 // k = L[R[p-1]]
for p = 0, ..., n-1
    k ← max(0, k-1) // safe dec.
    i ← R[p]
    q ← S[i+1]
    while (p+k < n) & (q+k < n) & (s[p+k] == t[q+k])
        k ← k+1 // inc.
    L[i] ← k
    
```

increments = # decrements + final value - initial $\leq 2n \Rightarrow$ runs in $O(n)$ time

Construction of the suffix array (by doubling) DSM/5

idea: compare suffixes by prefixes of length 2^i , for $i=0, \dots, \lceil \log n \rceil$

Def. $p =_k \delta$ if $p[1:k] = \delta[1:k]$ (by prefixes of length k)

$p \leq_k \delta$ if $p[1:k] \leq \delta[1:k]$
 \uparrow lex.

Note: \leq_k is only quasi-order (not necessarily antisymmetric),

we may have $p =_k \delta$ although $p \neq \delta$, e.g. $abc =_2 abd$.

Def. $R_k[i] = |\{j \mid \alpha[j:] \leq_k \alpha[i:]\}|$, i.e. # strictly smaller (w.r.t. \leq_k) suffixes than $\alpha[i:]$.

(“approximation of rank array”, $R_n = R$) first half

$\alpha[i:] \leq_{2k} \alpha[j:] \iff \alpha[i:] \leq_k \alpha[j:]$ on ($\alpha[i:] =_k \alpha[j:]$ and $\alpha[i+k:] \leq \alpha[j+k:]$) second half

$R_{2k}[i] \leq R_{2k}[j] \iff R_k[i] < R_k[j]$ or ($R_k[i] = R_k[j]$ and $R_k[i+k] \leq R_k[j+k]$)

$\iff (R_k[i], R_k[i+k]) \leq (R_k[j], R_k[j+k])$

\uparrow lexicographically as a pair

Alg. $k=1$: compute R_1 by sorting of suffixes in the first symbol
 $\implies O(n \log n)$ time (if no assumption on alphabet)

$k \rightarrow 2k$: sort suffixes by pairs $(R_k[i], R_k[i+k])$

+ S from R **Bucket sort**: two passes with $n+1$ buckets
in $O(n)$ time

$\implies O(n)$ time and space per step

$\lceil \log n \rceil$ steps $\implies O(n \log n)$ total time, $O(n)$ space

Note: It is possible to construct suffix array in $O(n)$ time and $O(1)$ additional space. (omitted)