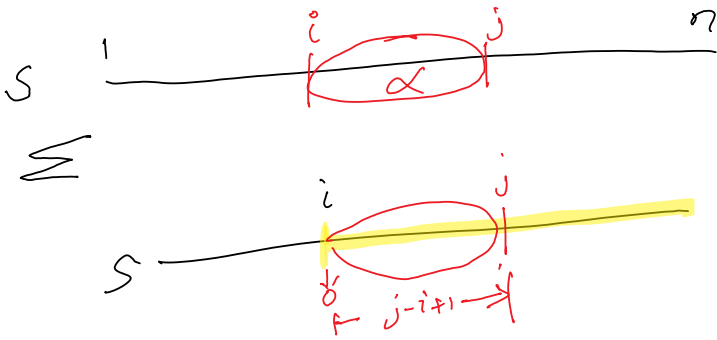


Suffix Tree Applications

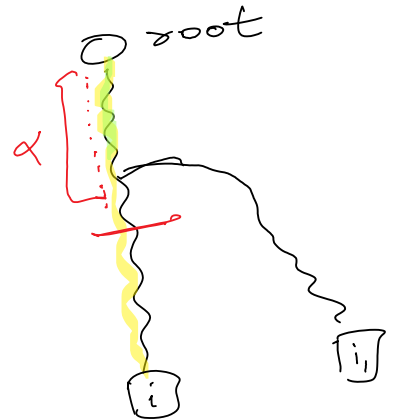
Monday, March 15, 2021 10:42 AM



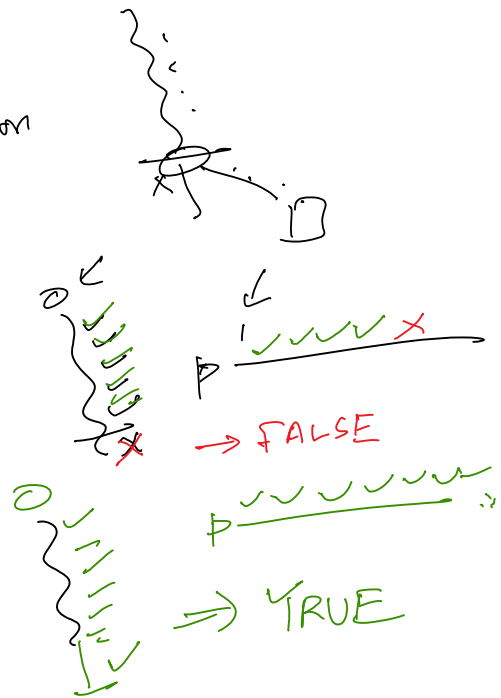
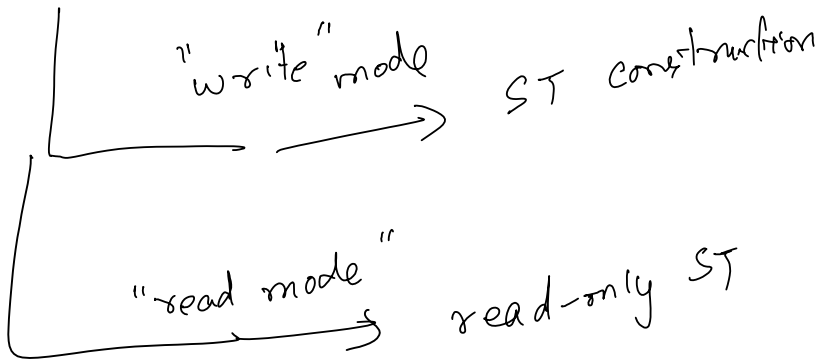
Every substring α of S , will have a unique path from the root \Rightarrow "index".

Suffix Trees

- 1) Compacted trie of all suffixes in S
- 2) String index for all substrings in S .



Find Path



Pattern Matching

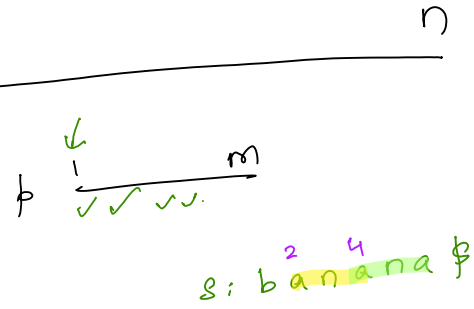
Monday, March 15, 2021 11:01 AM

(practical assumption: $n \gg m$)

Problem: Given a string $s[1..n]$ and a query pattern $p[1..m]$, find if pattern p occurs as a substring in s .

Input: $s[1..n], p[1..m]$

Output: i) Yes or No (decision problem) \rightarrow yes
ii) starting indices (location) $\rightarrow \{2, 4\}$

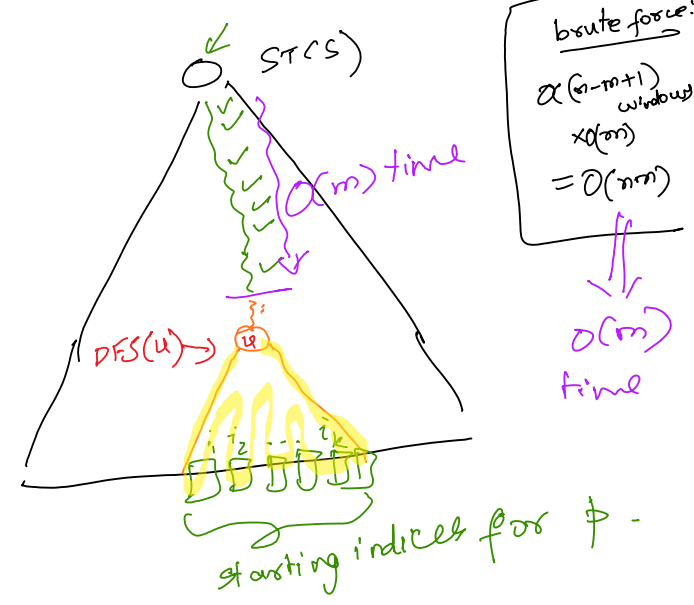


Algorithm:

1) Build STCS $\Rightarrow O(n)$ time
2) Find p in $s \Rightarrow O(m)$ time

FindPath(root, $p[1..m]$)
"read-only" mode

if p is exhausted \Rightarrow return FOUND
else // mismatch \Rightarrow return NOTFOUND



$O(n + m)$
preprocessing cost \uparrow querying cost

3) Go down to the next internal node (u) \Rightarrow size of the subtree under u . $\Rightarrow O(k)$ time
& enumerate all leaves under u . \Rightarrow output size

Q) Try to prove that this algorithm is correct: a) if p is in s , \Rightarrow algo will find it.
b) if algo outputs Yes/No $\Rightarrow p$ is in s (or not)

Longest Repeat

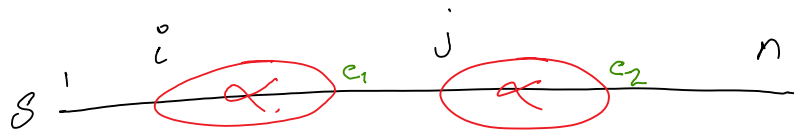
Monday, March 15, 2021 10:44 AM

Problem: Given a string s of length n (over alphabet Σ), find a longest repeat in s .

Definition: A "repeat" in a string s is a substring that occurs at least twice in s .

Input: $s \in \{1..n\}$

Output: longest Repeat α .

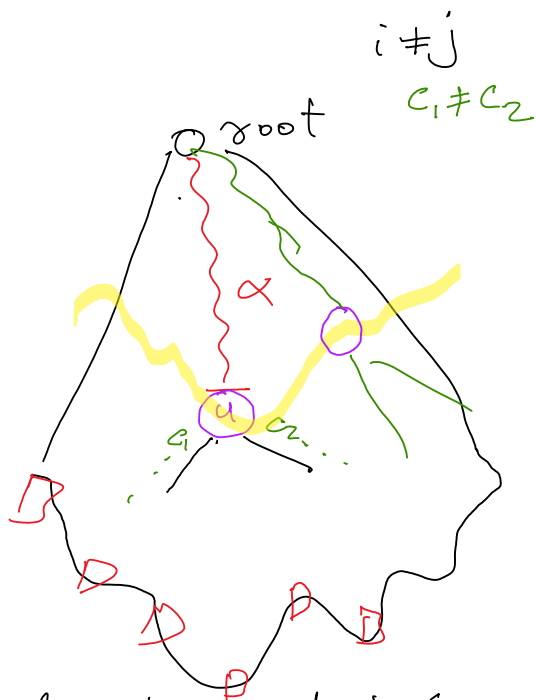


Algorithm:

1) Build $ST(s)$

2) a) Locate the "deepest" internal node (u) by string-depth

b) Output path-label (u)



Proof of correctness:

If algorithm outputs α

$\Rightarrow \alpha$ is a pathlabel of u

$\Rightarrow \alpha$ is a repeat $\rightarrow (1)$

Since u is the deepest

$\Rightarrow \alpha$ is a longest repeat

\square

If α is a longest repeat in s

\Rightarrow it is present in at least 2 places $i \& j$ s.t. $c_1 \neq c_2$

\Rightarrow There has to exist an internal node u at the end of the path labelled α

$\Rightarrow u$ is also the deepest (since α is the longest)

\Rightarrow Algorithm will output α .

\square

Suffix-Prefix Match

Monday, March 15, 2021 10:50 AM

A "superstring" is a string that contains others strings.

Problem:

Shortest Superstring Problem \Rightarrow Suffix-Prefix Exact Matching

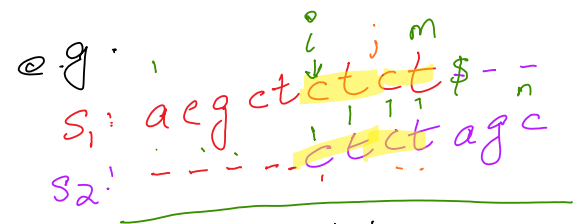
Algo: $O(m+n)$ time

Input: $s_1[1..m], s_2[1..n]$

Output: Find the longest suffix-prefix matching between s_1 & s_2 .

Algo: // first compute superstring for case (s_1 on left)

- 1) Build ST (s_1)
- 2) Do FindPath (root, $s_2[1..n]$) going as deep as possible until first mismatch
- 3) If matching stopped at a $\$$ -leaf
 - let $i \leftarrow$ label of the $\$$ -leaf
 - superstring = $s_1[1..m] \cdot s_2[m-i+2..n]$



\Rightarrow shortest superstring = $a e g c t c t c t a g c$

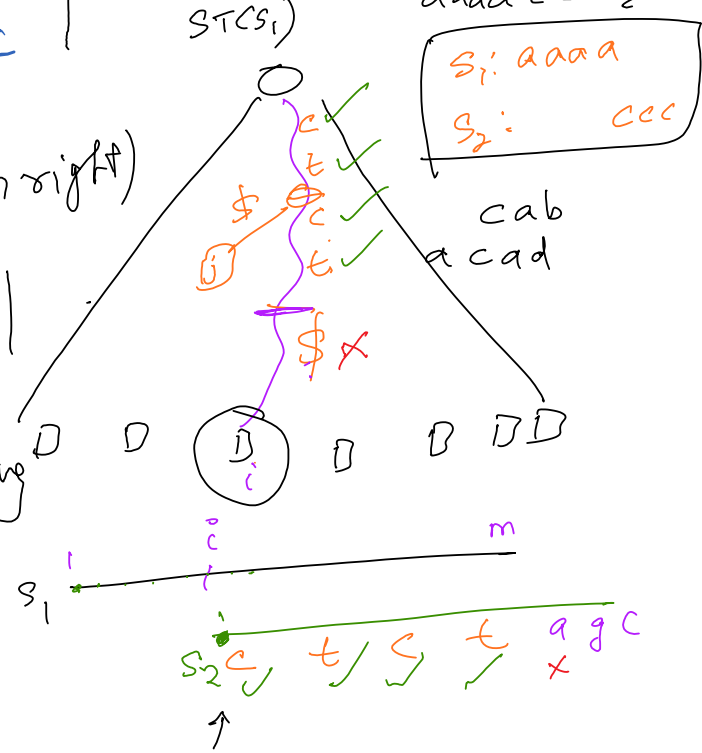
\Rightarrow output $\langle s_1, i \rangle$

// Next compute for case s_2 on left

- 4) Build ST (s_2)
- 5) Repeat steps 2 & 3 for (s_1 on right)
- 6) $|superstring_{1,2}|$ vs. $|superstring_{2,1}|$

Shorter of the two \Rightarrow superstring

if no superstring found \Rightarrow Not found

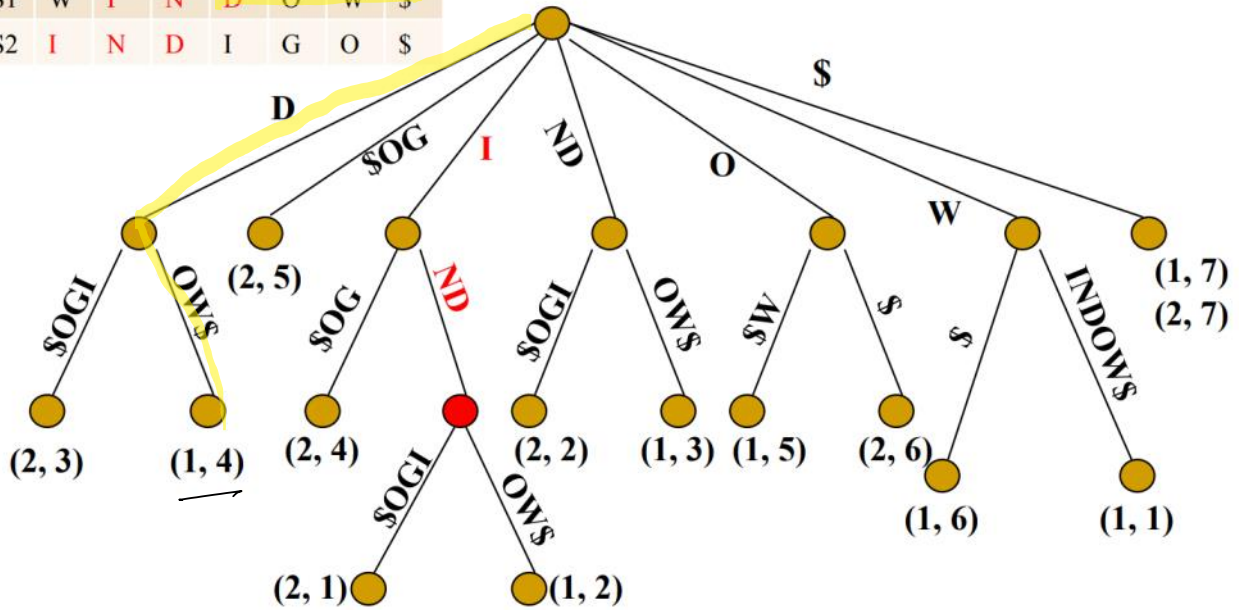


Generalized Suffix Trees (GST)

Monday, March 15, 2021 10:44 AM

Definition: Given a set $S = \{s_1, s_2, \dots, s_k\}$ of k strings, the Generalized Suffix Tree (GST) of S is the compacted trie of all suffixes of all strings in S .

	1	2	3	4	5	6	7
S1	W	I	N	D	O	W	\$
S2	I	N	D	I	G	O	\$



Leaf label format:
 <string id, suffix id>

Longest Common Substring (LCS)

Monday, March 15, 2021 10:45 AM

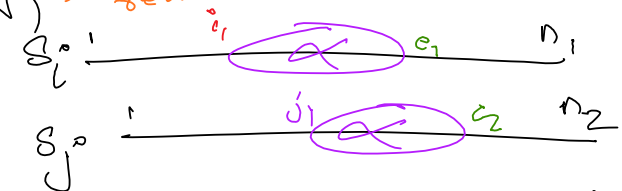
Problem: Given K strings, find an lcs between every pair of strings.

Input: $\mathcal{S} = \{s_1, s_2, \dots, s_k\}$ of lengths $\{n_1, n_2, \dots, n_k\}$ resp.

Output: For each pair $\langle s_i, s_j \rangle$ find lcs. $N = \sum_{i=1}^k |s_i|$

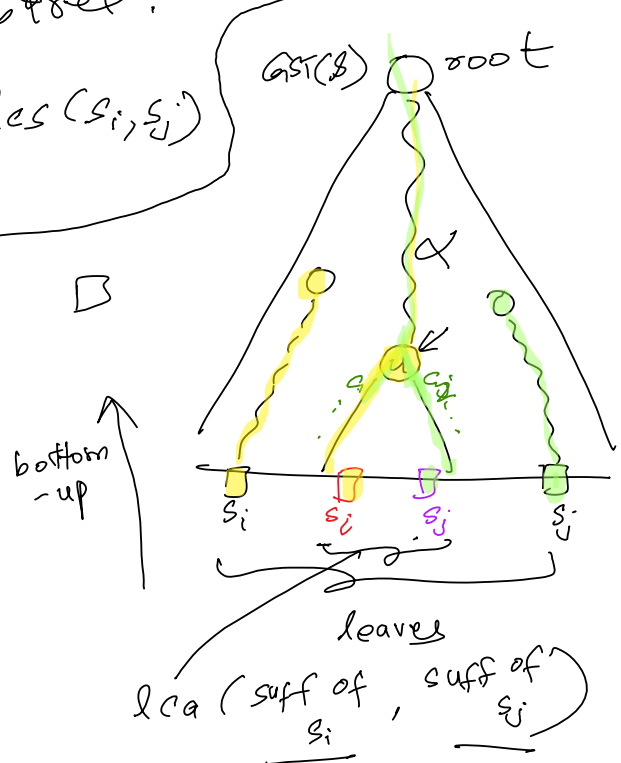
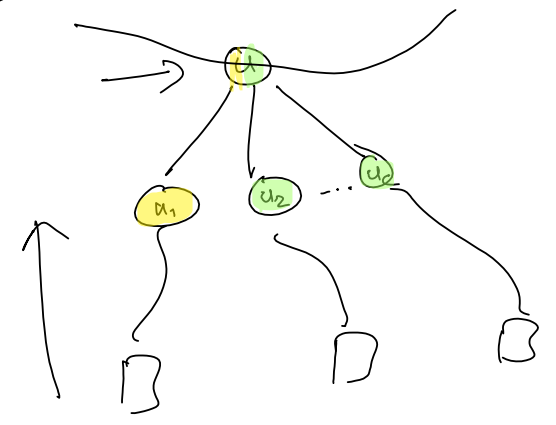
Algo:

- 1) Build $\text{GST}(\{s_1, s_2, \dots, s_k\}) \Rightarrow O(N)$ time
- 2) For each $\langle s_i, s_j \rangle \in \mathcal{S}$ $\Rightarrow O(N)$ → can it be reduced?



a) find the deepest internal node u that has at least one leaf from s_i & one leaf from s_j under its subtree.

b) Report $\text{pathlabel}(u) \rightarrow \text{lcs}(s_i, s_j)$



LCP problem revisited, using LCA

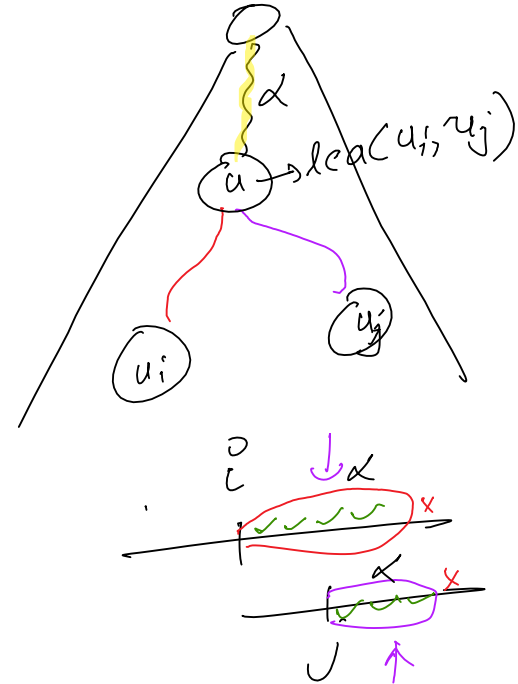
Friday, March 19, 2021 10:54 AM

Lowest Common Ancestor
 $LCA(u_i, u_j) \Rightarrow lca$

Function

$lcp \equiv lca$
↓
longest common prefix
($suff_i, suff_j$)

↓
lowest common ancestor
(u_i, u_j) in the tree



String \Rightarrow reduces to
tree
 $lca = ? \quad O(i)$

Suffix-Prefix Match (revisited) → but using GSTs

Monday, March 15, 2021 10:50 AM

Problem:

Shortest Superstring Problem ⇒ Suffix-Prefix Exact Matching

Input: $s_1 [1..m], s_2 [1..n]$

Output: Find the longest suffix-prefix matching



Alg^o:

- Build GST(s_1, s_2)

2.

$l \equiv \text{lcp}(\text{suff}_i \text{ of } s_2, \text{suff}_i \text{ of } s_1)$

if $(l == m - i + 1)$

Yes

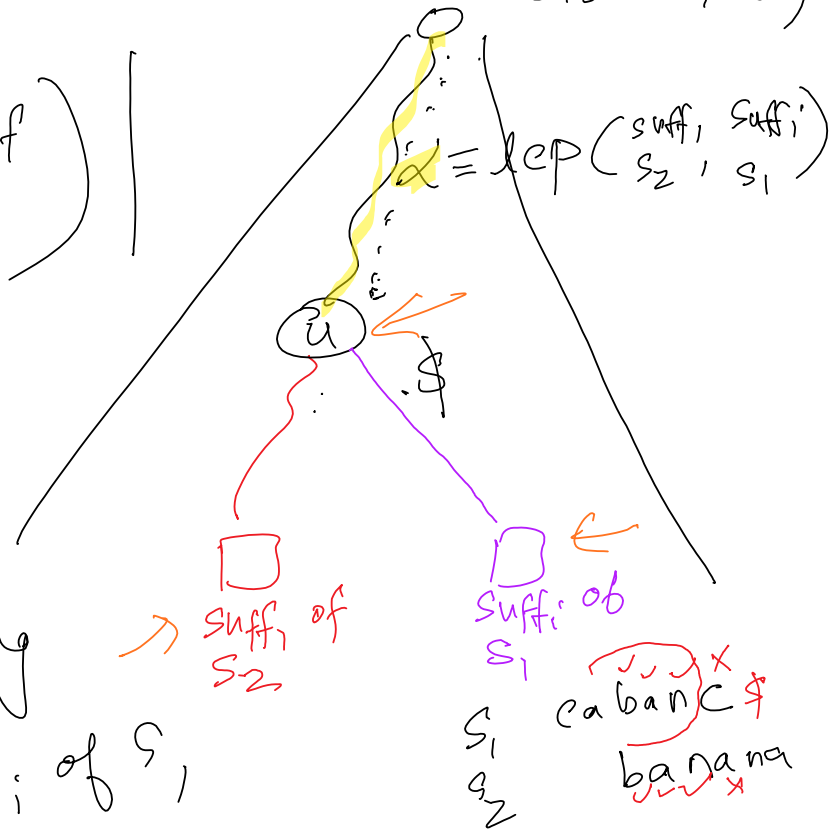
superstring exists

$s_1 \rightarrow s_2$

NO

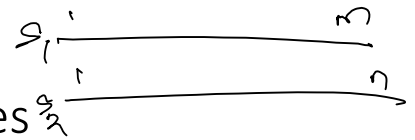
no superstring exist from suff_i of s_1

involves suff_i



Suffix-Prefix Match - revisited using LCA queries

Monday, March 22, 2021 10:52 AM



$$\text{lep}(\text{suff}_i, \text{suff}_j) \equiv \text{lca}(\text{leaf}_i, \text{leaf}_j)$$

ALGORITHM:

Build GST (s_1, s_2)

// check if there exists a superstring with $(s_1 \text{ on left})$

for $(i = 1 \text{ to } m)$

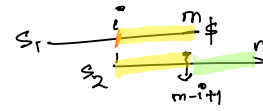
$l \leftarrow \text{lep}(\text{suff}_i \text{ of } s_1, \text{suff}_1 \text{ of } s_2)$

check if $(l == m - i + 1)$

YES: $\text{super}_{1,2}$
break;

NO: continue;

If $\text{super}_{1,2}$ not found inside the loop,
then $\text{super}_{1,2} \leftarrow s_1[1..m] \cdot s_2[1..n]$ // set trivial



// check if there exists a superstring with $(s_2 \text{ on the left})$

for $(j = 1 \text{ to } n)$

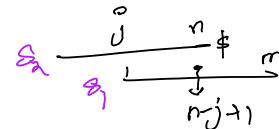
$l \leftarrow \text{lep}(\text{suff}_j \text{ of } s_2, \text{suff}_1 \text{ of } s_1)$

check if $(l == n - j + 1)$

YES: $\text{super}_{2,1}$
break;

NO: continue;

If $\text{super}_{2,1}$ not found inside the loop,
then $\text{super}_{2,1} \leftarrow s_2[1..n] \cdot s_1[1..m]$ // set trivial



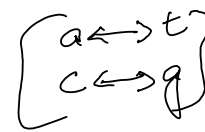
// Report the shortest

superstring $\leftarrow \text{shortest} \{ \text{super}_{1,2}, \text{super}_{2,1} \}$

B.

Restriction Enzymes

Monday, March 15, 2021 10:50 AM



cut site

palindrome

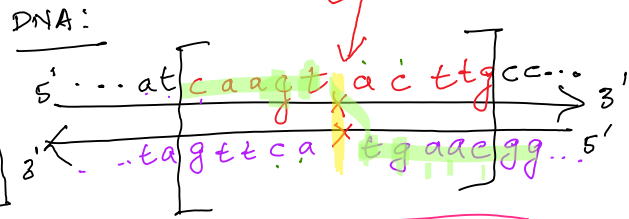
← snadaM →

← racecar →

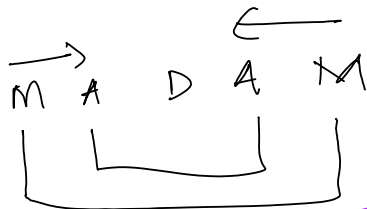
← matajalam →

← ABBA →

Restriction enzymes cut at sites that have DNA palindromes.



Definition: A "DNA palindrome" is a sequence s of some length l , such that $s[i+j] = \text{rev.comp}(s[l-i-j])$, $\forall 0 \leq i \leq \lfloor \frac{l}{2} \rfloor$

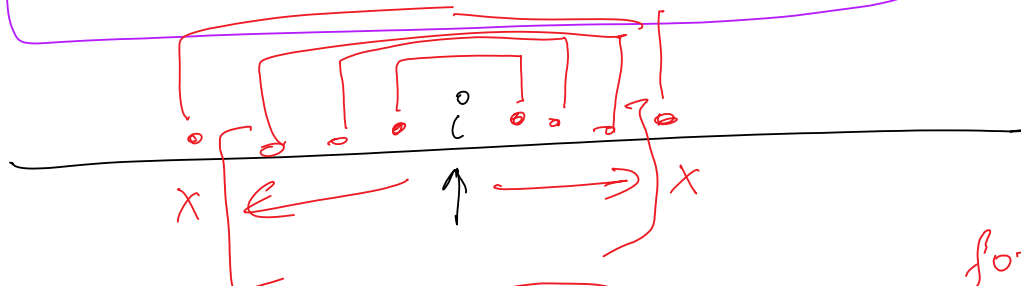


e.g. $s = \text{caagctacttg}$

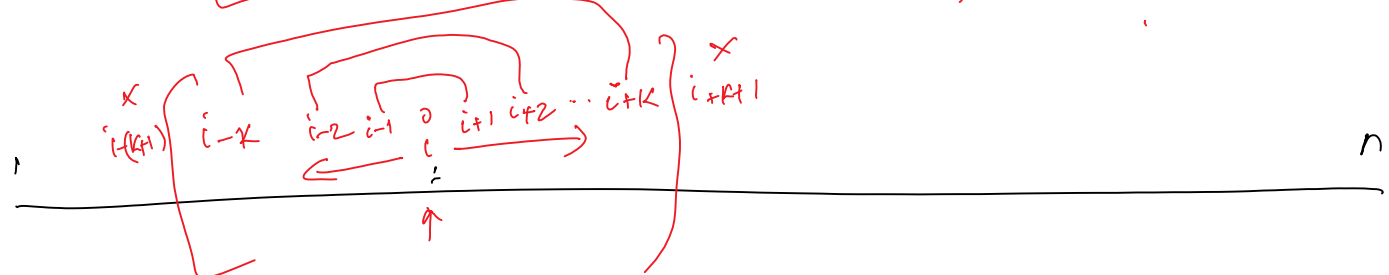
e.g. $s = \text{caagccttg}$

I have a racecar that is fast.

Build $\text{GST}(s, s^r)$



for $i = 1$ to n



RE Site Detection: Algorithm

Wednesday, March 24, 2021 10:51 AM

Input: A genomic sequence $S[1..n]$.
 Output: All palindromic sites of length $\geq \tau_{min}$

Algorithm:

1) Build $GST(S, S^r)$ $\Rightarrow O(n)$ time

// odd palindrome case. $\Rightarrow O(n)$ time

2) for ($i = 1$ to n) { $\Rightarrow O(i)$ time

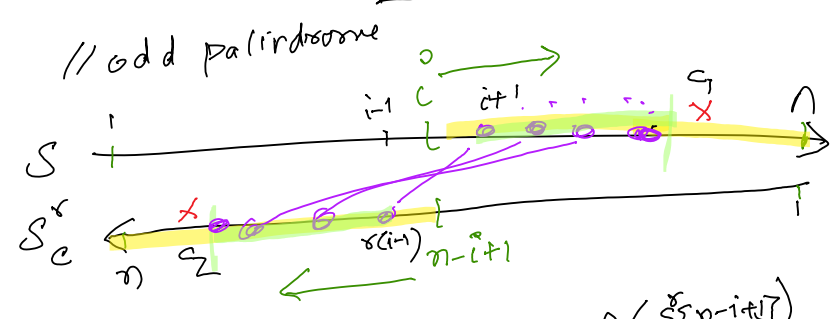
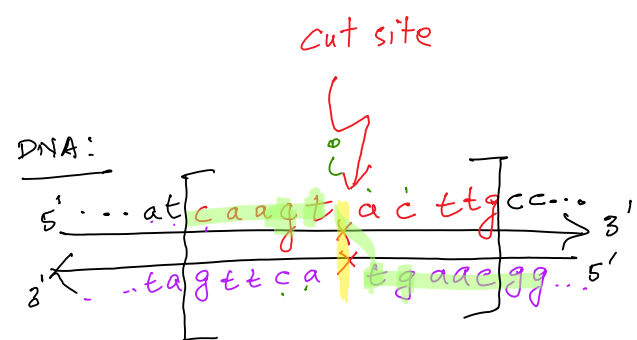
a. $u \leftarrow \text{lca}(\langle S, i+1 \rangle, \langle S^r, \tau(i-1) \rangle)$

b. $lcp \leftarrow \text{pathlabel}(u)$

c. if ($|lcp| \geq \tau_{min}$) {
 print "odd palindrome centered at i exists"
 palindrome = $S[i+1..i+|lcp|-1]$

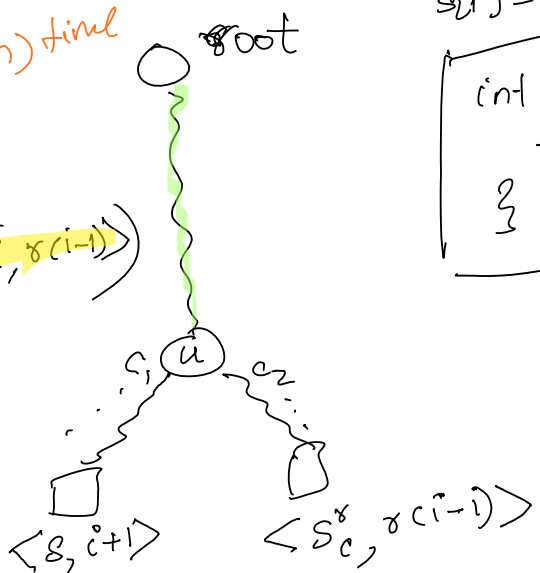
// do similar check for even length palindromic case.

$\Rightarrow O(n)$ time total



```

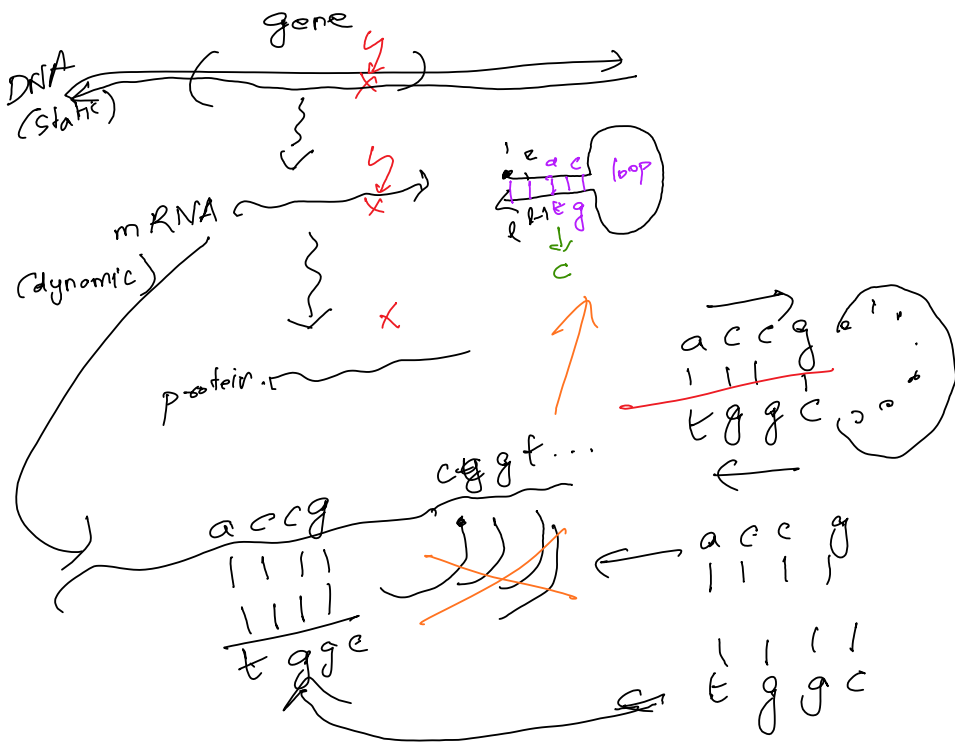
S[i] = comp (S[n-i+1])
int r(i) {
    return (n-i+1)
}
    
```



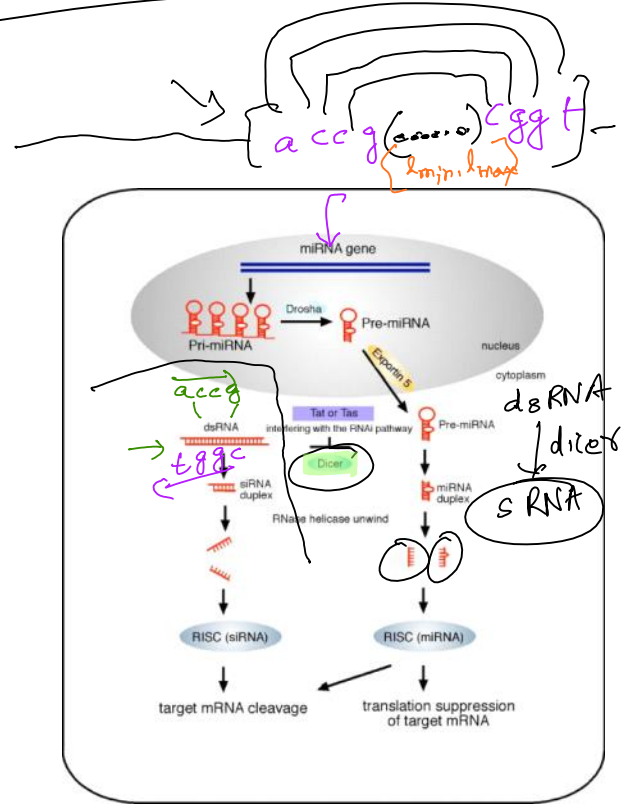
Gene Silencing

Wednesday, March 17, 2021 10:53 AM

Q) How to "silence" a gene?



RNA interference: Nobel prize, 2006 (Fire & Mello)



DOI: [10.1186/1742-4690-2-35](https://doi.org/10.1186/1742-4690-2-35)

