

Tries

Friday, February 26, 2021 10:50 AM

Motivation: To capture variable-length matches (directly).

$$\mathcal{S} = \{s_1, s_2, s_3, \dots, s_k\}$$

	1	2	3	4	5	6
s_1 :	a	p	p	l	e	
s_2 :	a	p	p	e	a	r
s_3 :	a	p	p	l	y	
s_4 :	a	p	p			

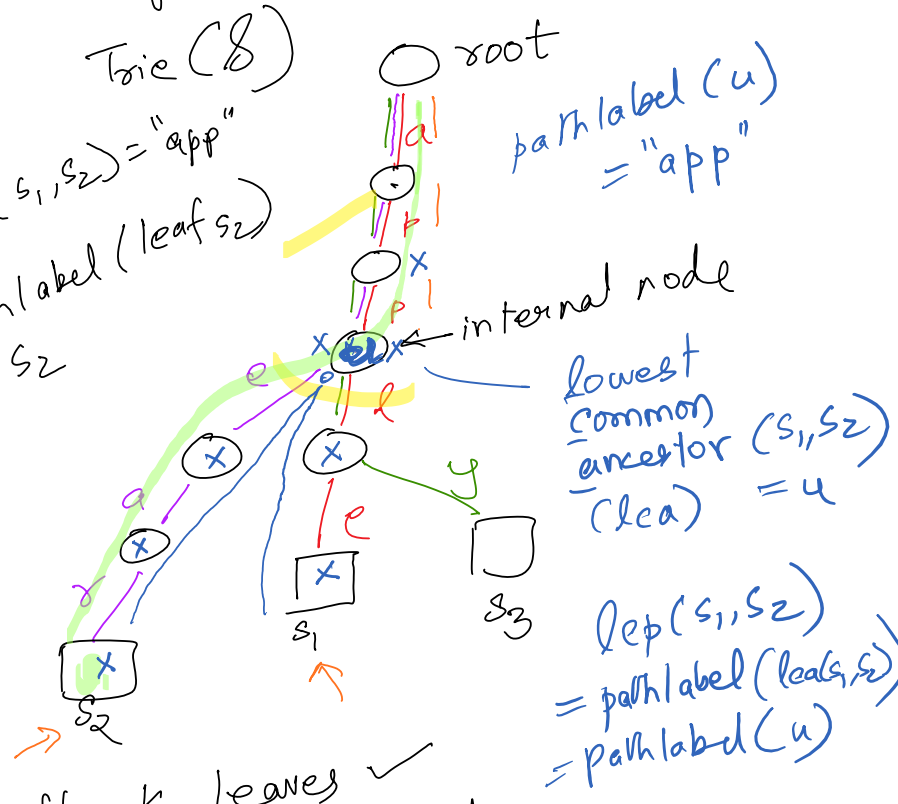
Definition:

Given $\mathcal{S} = \{s_1, \dots, s_k\}$ from alphabet Σ , a trie for \mathcal{S} is:

- a) rooted tree with exactly k leaves ✓
- b) each edge is labeled with a (single) character.
- c) no two outgoing edges of an internal node start with the same character label
- d) the path label from the root to any leaf should spell out one of the input strings.

Trie(\mathcal{S})

$lep(s_1, s_2) = \text{"app"}$
 $pathlabel(\text{leaf } s_2) = s_2$



For a compacted trie, this rule is relaxed so edge labels can be strings (of length ≥ 1).

Appending \$ symbol to make the leaves unique

Monday, March 1, 2021 11:14 AM

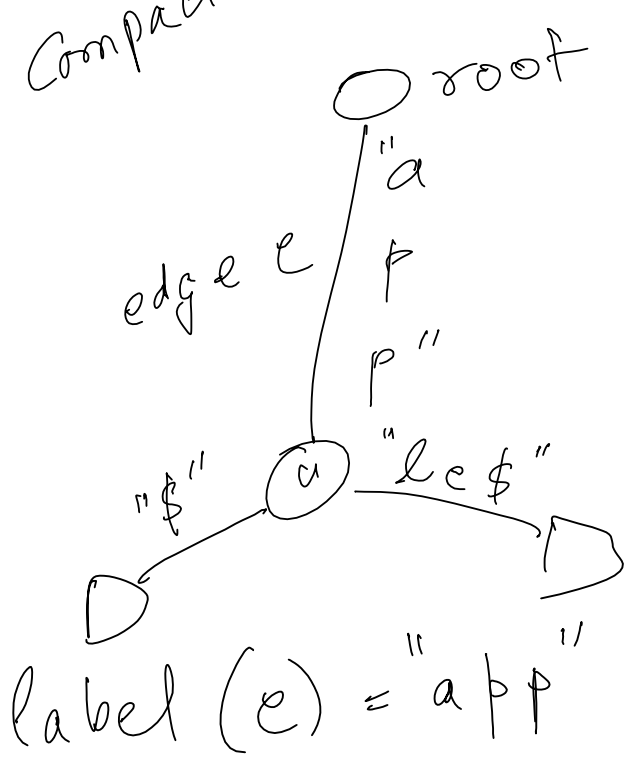
$$\Sigma = \{a, c, g, t\}$$

Add a new symbol $\$$ $\notin \Sigma$ the input set S .
 Append $\$$ to all strings in S .

$$S = \{s_1, s_2, s_3, \dots, s_k\}$$

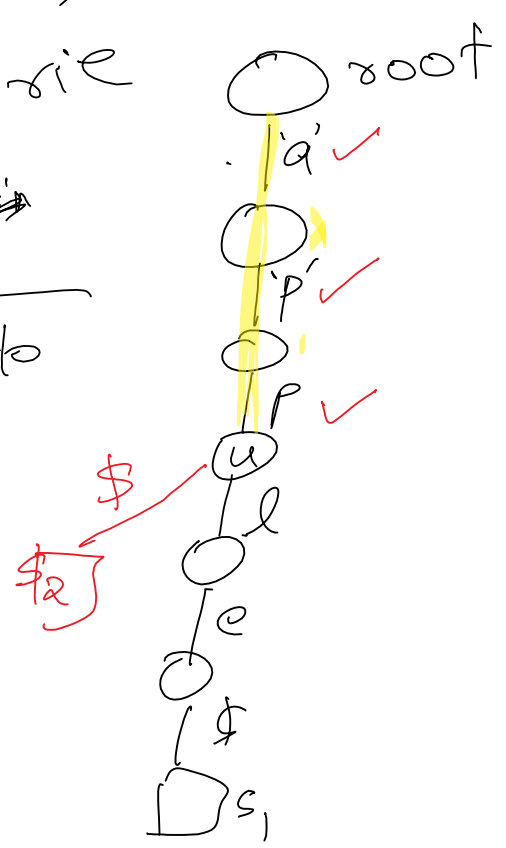
$\rightarrow s_1 = \text{apple}\$$
 $\rightarrow s_2 = \text{app}\$$

Compacted trie



Trie

compactification
 edges to
 eliminate
 int. nodes
 with only
 one child



Practical Algorithms for Text Retrieval for Info. Coded in Alphanumeric.

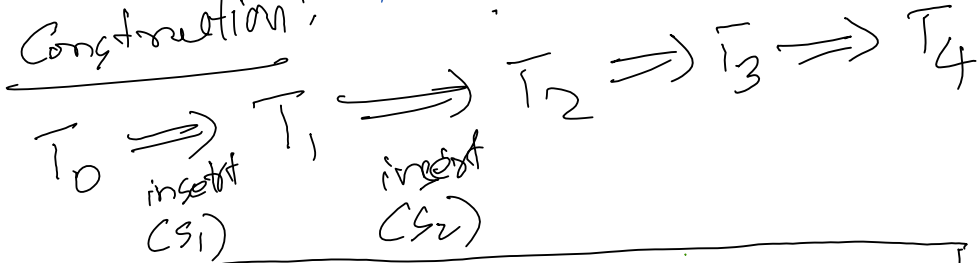
PATRICIA trees

Friday, February 26, 2021 10:50 AM

PATRICIA tree: = the compacted trie for a set of strings S .

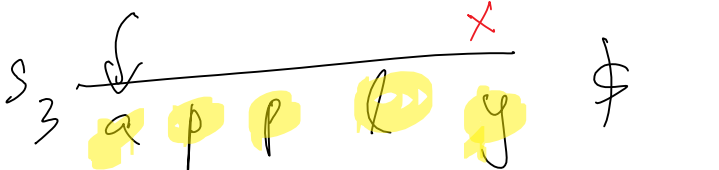
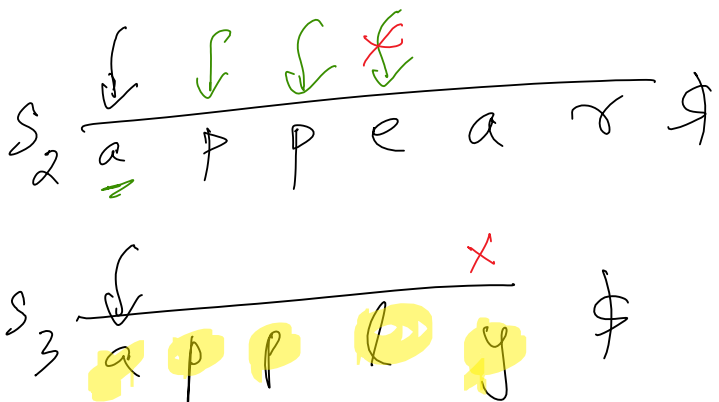
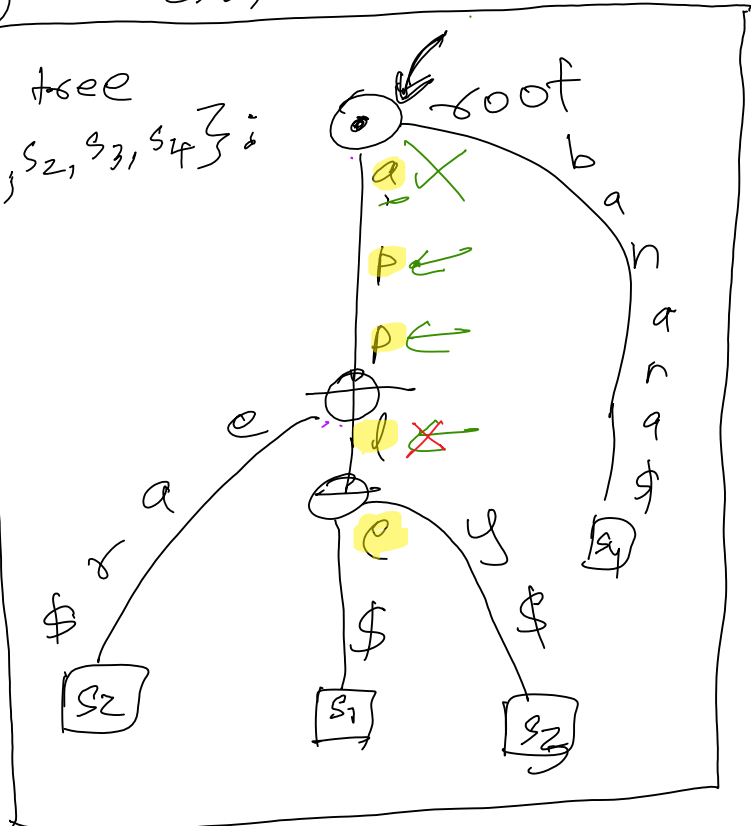
Input: $S = \{s_1, s_2, \dots, s_k\}$

Construction: // insert one string at a time

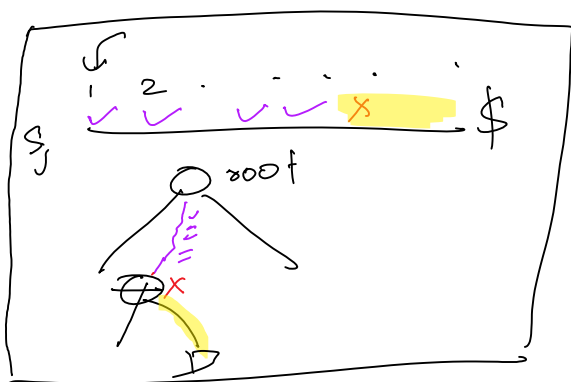


	1	2	3	4	5	6
s_1 :	a	p	p	l	e	
s_2 :	a	p	p	e	a	r
s_3 :	a	p	p	l	y	
s_4 :	b	a	n	a	n	a

PATRICIA tree for $\{s_1, s_2, s_3, s_4\}$:



Insert of s_j into T
 Find Path (root, $\langle s_j, 1 \rangle$)

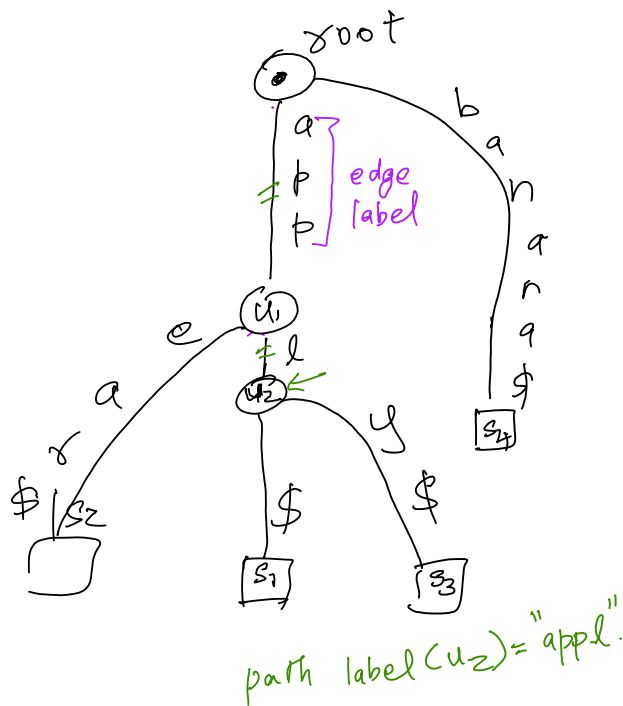


Trie Definitions

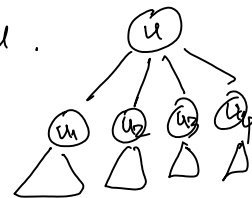
Wednesday, March 3, 2021 10:42 AM

string-depth(u_2) = 4
node-depth(u_2) = 2

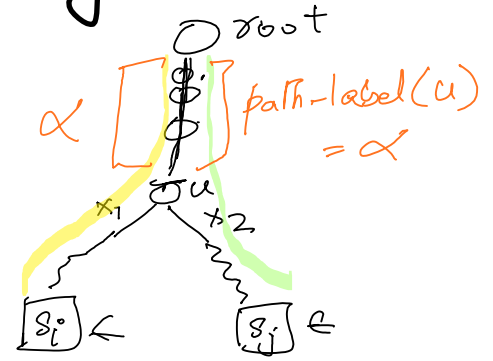
- Node: either internal node or leaf
- Internal node: node with ≥ 2 children
- Leaf: node with \emptyset children
- Root: an internal node (special case w/ 0 or more children)



- Edge-label: string label of an edge
- Path-label(u): the concatenation of all edge labels from the root to u .
- String-depth(u): length of path-label(u)
- Node-depth(u): Number of edges between the root and u .
- Parent(u): The parent node of node u (if it exists) (root has no parent)
- Ancestor(u): Any node on the path from u to the root.
- child(u): Any (immediate) child node of u .



Q1) What is the longest common prefix (lcp) between any two strings? $\langle s_i, s_j \rangle$



Let α be the lcp between s_i & s_j .

\Rightarrow Then there has to exist an internal node u , with path-label α under which, $\text{leaf}(s_i)$ and $\text{leaf}(s_j)$ are present in two different branches (as shown)

\Rightarrow In other words, look for the lowest common ancestor (lca) of the $\text{leaf}(s_i)$ and $\text{leaf}(s_j)$ and output its path-label.

Input: $S = \{s_1, s_2, \dots, s_i, \dots, s_j, \dots, s_k\}$

Algorithm:

- 1) Build the PATRICIA tree for S // one time preprocessing
- 2) Query time: // querying for lots of lcp by user
 - a) user provides $\langle s_i, s_j \rangle$ pair
 - b) Compute $\text{lca}(\text{leaf}(s_i), \text{leaf}(s_j)) \rightarrow \text{node } u$
 - c) Output path-label(u).

Longest Common Prefix ==> Lowest Common Ancestor

Monday, February 24, 2020 10:53 AM

LCP problem:

Input: $S = \{s_1, s_2, \dots, s_k\}$

Output: $lcp(s_i, s_j)$ query

Alg 0: \rightarrow preprocessing step
1. Build PATRICIA (S)

2. locate s_i & s_j leaves

color $s_i \leftarrow$ purple

color $s_j \leftarrow$ red

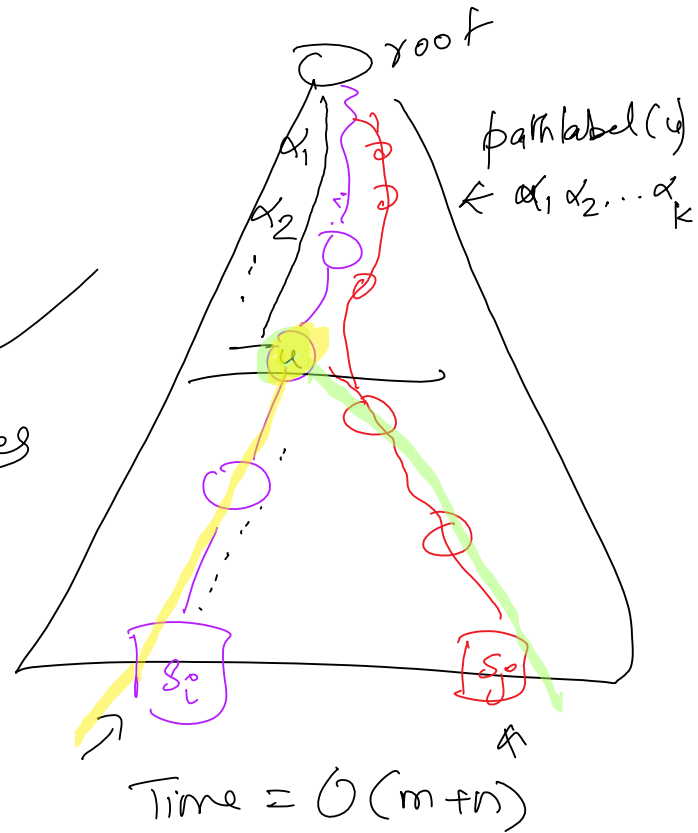
color all ancestors of s_i
w/ purple

" " " s_j
w/ red

3. $lea(s_i, s_j) \leftarrow$ lowest / deepest node which
has both colors ("mixed node")

Let $u \leftarrow lea(s_i, s_j)$

4. Output $pathlabel(u)$.



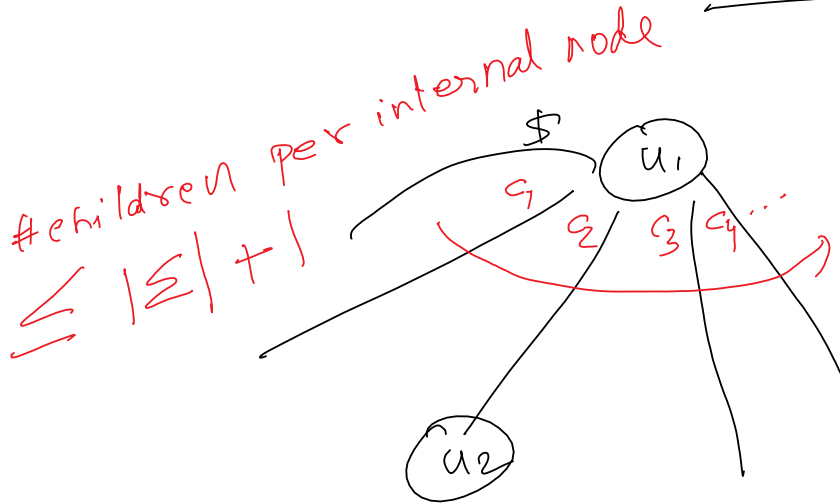
PATRICIA : $\Sigma = \{s_1, s_2, \dots, s_k\}$

Trie properties & implementation

Monday, March 1, 2021 11:45 AM

$\Sigma \neq \emptyset$

Compacted trie (or PATRICIA tree) is a $(|\Sigma|+1)$ -way tree.



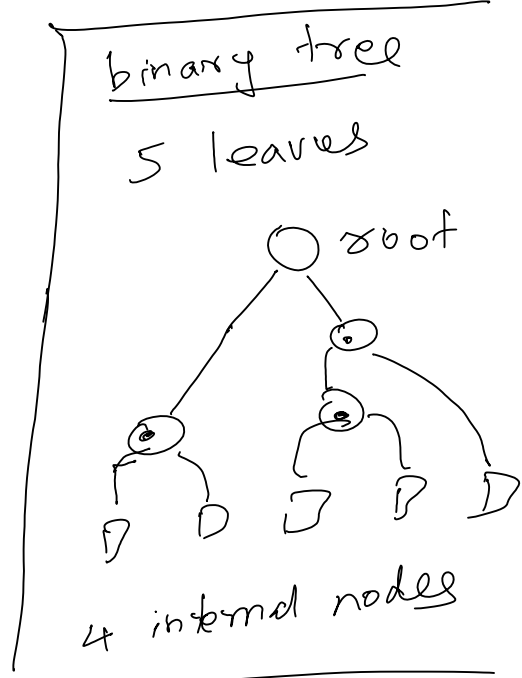
Tree:

- root
- internal nodes
- leaves
- edges

edge labels
(strings over Σ)

$\Sigma = \{a, c, g, t\}$

- 1) 1 root (special case of an internal node)
- 2) Number of leaves = k
- 3) Number of internal nodes $\leq k - 1$

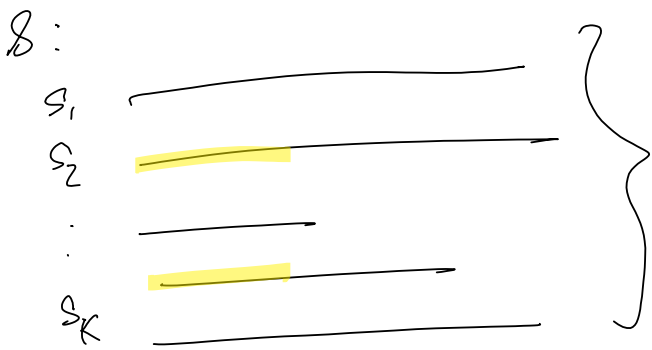


Exact matching use-cases: LCP, LCS (and a reason to build a trie for suffixes)

Wednesday, March 3, 2021 10:55 AM

(and a reason to build a trie for suffixes)

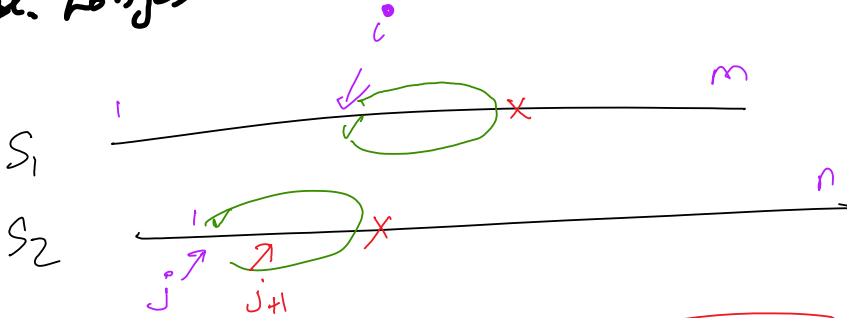
Use case: Longest Common Prefix (LCP)



$$\text{lcp}(s_i, s_j) = ?$$

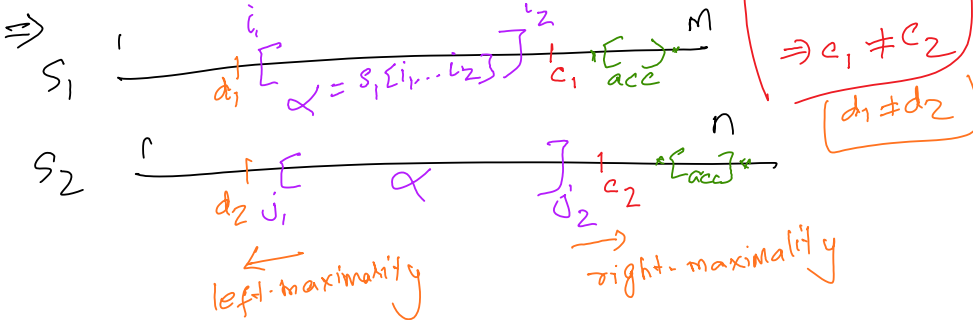
- 1) Build PATRICIA tree
- 2) $\text{lca}(s_i, s_j) \rightarrow u$
output $\text{pathlabel}(u)$

Use case: Longest Common Substring (LCS)

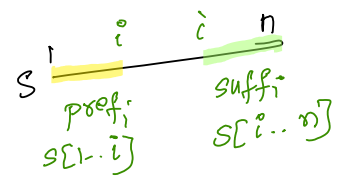


Challenge:
But PATRICIA trees can only help capture matches that are prefixes.

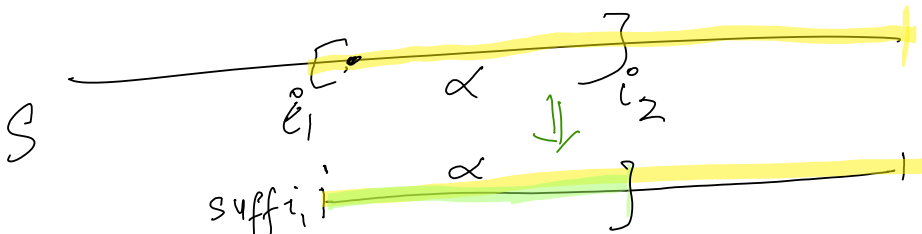
Let $\alpha = \text{lcs}(s_1, s_2)$



$\text{lcs} \Rightarrow$ right & left maximal.
 \Rightarrow "maximal"



Q) What is a substring? \equiv a prefix of a suffix



\Rightarrow so let us build a PATRICIA tree for all suffixes!