

MAXIME CROCHEMORE

King's College London

Maxime.Crochemore@kcl.ac.uk
<http://www.dcs.kcl.ac.uk/staff/mac/>

Implementation of indexes



Implementation with efficient data structures

- ★ **Suffix Trees**
digital trees, PATRICIA tree (compact trees)
- ★ **Suffix Automata or DAWG's**
minimal automata, compact automata

Implementation with efficient algorithm

- ★ **Suffix Arrays**
binary search in the ordered list of suffixes

Suffixes

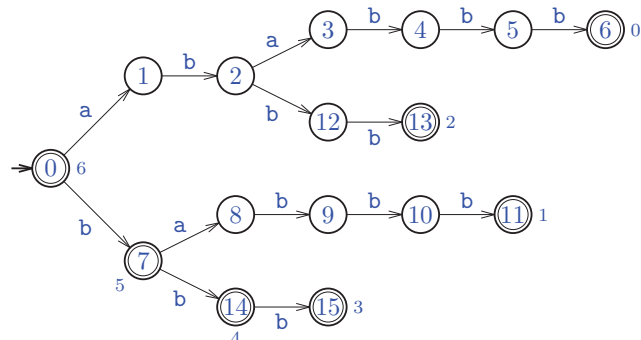
Text $y \in \Sigma^*$

- ★ $Suff(y)$ = set of suffixes of y ,
- ★ $\text{card } Suff(y) = |y| + 1$
- ★ $Suff(\text{ababbb})$

i	0	1	2	3	4	5	
$y[i]$	a	b	a	b	b	b	
							position
	a	b	a	b	b	b	0
		b	a	b	b	b	1
			a	b	b	b	2
				b	b	b	3
					b	b	4
						b	5
						ε	6 (empty string)

Trie of suffixes

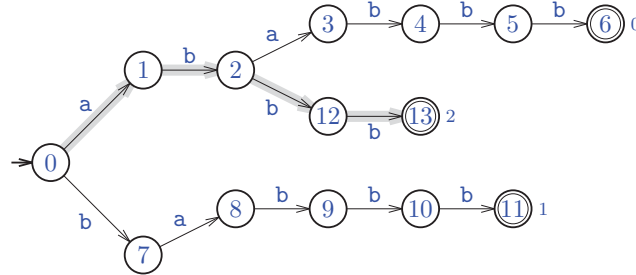
- ★ $\mathcal{T}(y)$ = digital tree which branches are labeled by suffixes of y
= tree-like deterministic automaton accepting $Suff(y)$
- ★ **Nodes**
identified with factors (subwords) of y
- ★ **Terminal nodes**
identified with suffixes of y , output = position of the suffix
- ★ **Suffix trie of ababbb**



Forks

Insertion of $u = y[i \dots n - 1]$ in the structure accepting longer suffixes

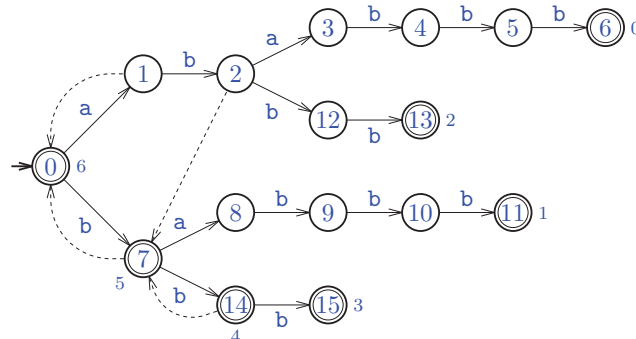
- ★ **Head** of u : longest prefix $y[i \dots k - 1]$ of u occurring before i
- ★ **Tail** of u : rest $y[k \dots n - 1]$ of suffix u
- ★ $y = \text{ababbb}$; head of **abbb** is **ab**; tail of **abbb** is **bb**



- ★ **Fork**
any node that has outdegree 2 at least,
or that both has outdegree 1 and is terminal
- ★ **Note**: the node associated with the head of u is a fork
initial node is a fork iff y non empty

Suffix link

- ★ **Function** s_y , *suffix link*
if node p identified with factor av , $a \in \Sigma$, $v \in \Sigma^*$
 $s_y(p) = q$, node identified with v



- ★ **Use**
creates shortcuts used to accelerate heads computations
- ★ **Useful for forks only**
undefined on initial node
- ★ **Note**: if p is a fork, so is $s_y(p)$

Suffix Tree

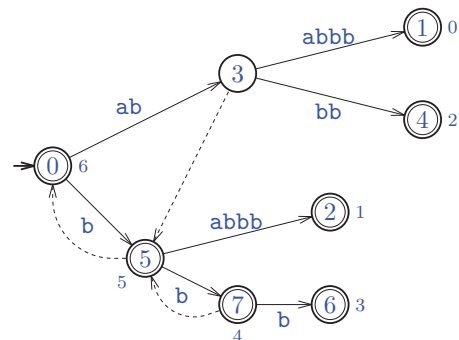
Text $y \in \Sigma^*$ of length n

$\mathcal{S}(y)$ suffix tree of y : compact trie accepting $\text{Suff}(y)$

★ **Definition**

tree obtained from the suffix trie of y by deleting all nodes having outdegree 1 that are not terminal

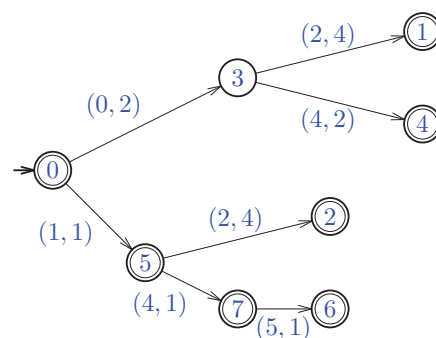
★ **Edges labeled by factors** of y instead of letters



- ★ **Number of nodes:** no more than $2n$ (if $n > 0$)
because all internal nodes have two children at least
and there are at most n external nodes

Labels of edges

★ **Labels represented by pairs** $(pos, Length)$



i	0	1	2	3	4	5
$y[i]$	a	b	a	b	b	b

- ★ Requires to have y in main memory
★ **Size of $\mathcal{S}(y)$:** $O(n)$

Scheme of suffix tree construction

SUFFIX-TREE(y)

```

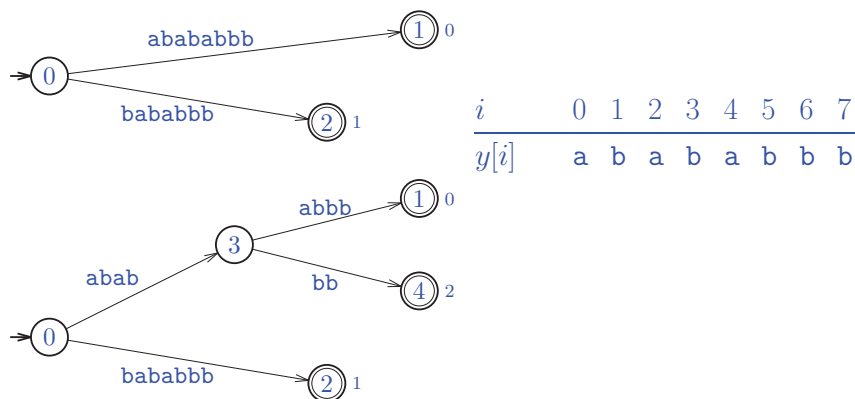
1   $T \leftarrow \text{NEW-TREE}()$ 
2  for  $i \leftarrow 0$  to  $n - 1$  do
3      find fork of head of  $y[i \dots n - 1]$  using
          FAST-FIND from node  $s[\text{parent}]$  if needed
          and then SLOW-FIND
4       $k \leftarrow$  position of tail of  $y[i \dots n - 1]$ 
5      if  $k < n$  then
6           $q \leftarrow \text{NEW-STATE}()$ 
7           $\text{Adj}[\text{fork}] \leftarrow \text{Adj}[\text{fork}] \cup \{(k, n - k), q)\}$ 
8           $\text{output}[q] \leftarrow i$ 
9      else  $\text{output}[\text{fork}] \leftarrow i$ 
10  $\text{output}[\text{initial}] \leftarrow n$ 
11 return  $T$ 

```

★ Adjacency-list representation of labeled arcs

Straight insertion

★ Insertion of suffix **ababbb** is done by letter comparisons from the initial node (current node)



★ It leads to create node 3 which suffix link is still undefined,
 ★ and node 4 associated with suffix **ababbb** at position 2
 ★ Head is **abab**, tail is **bb**

Slow find

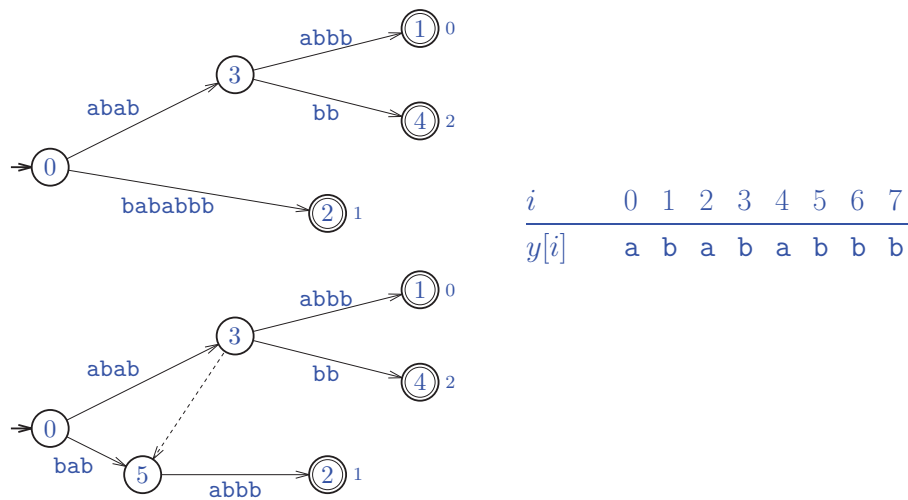
```

SLOW-FIND( $p, k$ )
1  while  $k < n$  and  $\text{TARGET}(p, y[k]) \neq \text{NIL}$  do
2       $q \leftarrow \text{TARGET}(p, y[k])$ 
3       $(j, \ell) \leftarrow \text{label}(p, q)$ 
4       $i \leftarrow j$ 
5      do  $i \leftarrow i + 1$ 
6           $k \leftarrow k + 1$ 
7      while  $i < j + \ell$  and  $k < n$  and  $y[i] = y[k]$ 
8      if  $i < j + \ell$  then
9           $\text{Adj}[p] \leftarrow \text{Adj}[p] \setminus \{((j, \ell), q)\}$ 
10          $r \leftarrow \text{NEW-STATE}()$ 
11          $\text{Adj}[p] \leftarrow \text{Adj}[p] \cup \{((j, i - j), r)\}$ 
12          $\text{Adj}[r] \leftarrow \text{Adj}[r] \cup \{((j + i - j, \ell - i + j), q)\}$ 
13         return  $(r, k)$ 
14      $p \leftarrow q$ 
15 return  $(p, k)$ 

```

New suffix link

- ★ Computing $s[3] = s_y(3)$ remains to find the node associated with **bab**



- ★ Arc $(0, (1, 7), 2)$ is split into $(0, (1, 3), 5)$ and $(5, (4, 4), 2)$
- ★ Execution in constant time (here)
- ★ In general, iteration in time proportional to the number of nodes along the path (and not proportional to the length of the string)

Fast find

FAST-FIND(r, j, k)

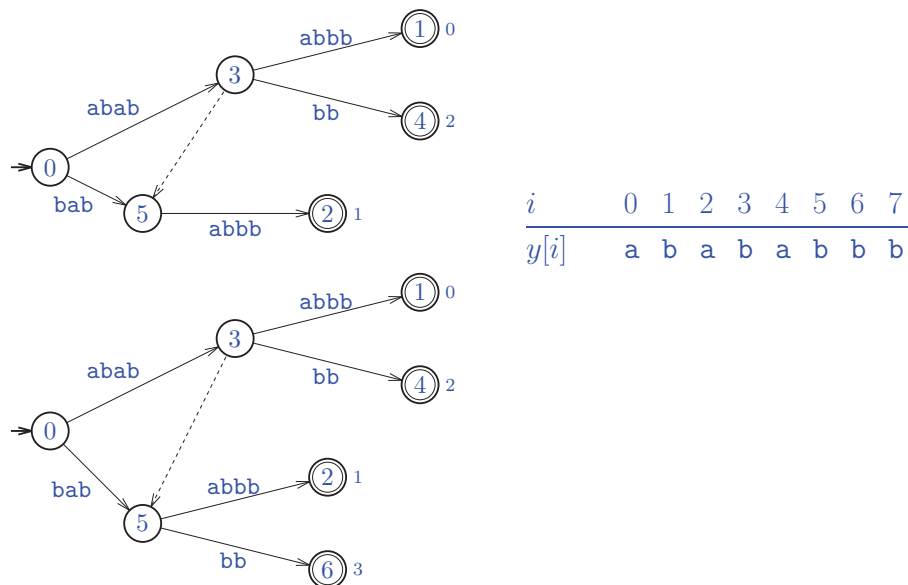
```

1  ▷ computes TARGET( $r, y[j \dots k - 1]$ )
2  if  $j \geq k$  then
3      return  $r$ 
4  else  $q \leftarrow$  TARGET( $r, y[j]$ )
5       $(j', \ell) \leftarrow$  label( $r, q$ )
6      if  $j + \ell \leq k$  then
7          return FAST-FIND( $q, j + \ell, k$ )
8      else  $\text{Adj}[r] \leftarrow \text{Adj}[r] \setminus \{((j', \ell), q)\}$ 
9           $p \leftarrow$  NEW-STATE()
10          $\text{Adj}[r] \leftarrow \text{Adj}[r] \cup \{((j, k - j), p)\}$ 
11          $\text{Adj}[p] \leftarrow \text{Adj}[p] \cup \{((j' + k - j, \ell - k + j), q)\}$ 
12         return  $p$ 

```

Next insertion

★ End of insertion of suffix **babbb**

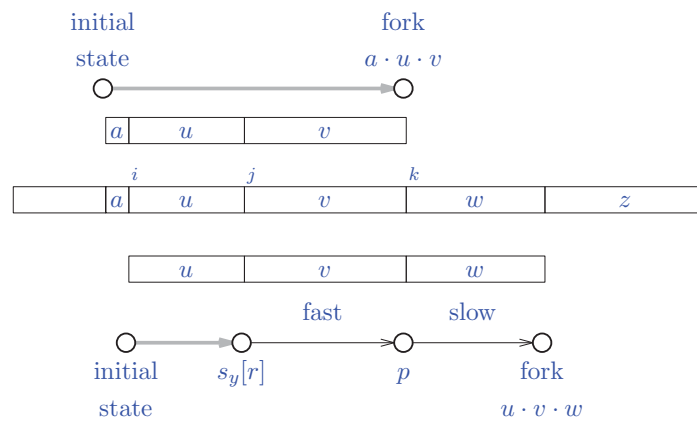


★ Execution in constant time

★ Head is **bab**, tail is **bb**

Scheme for insertion

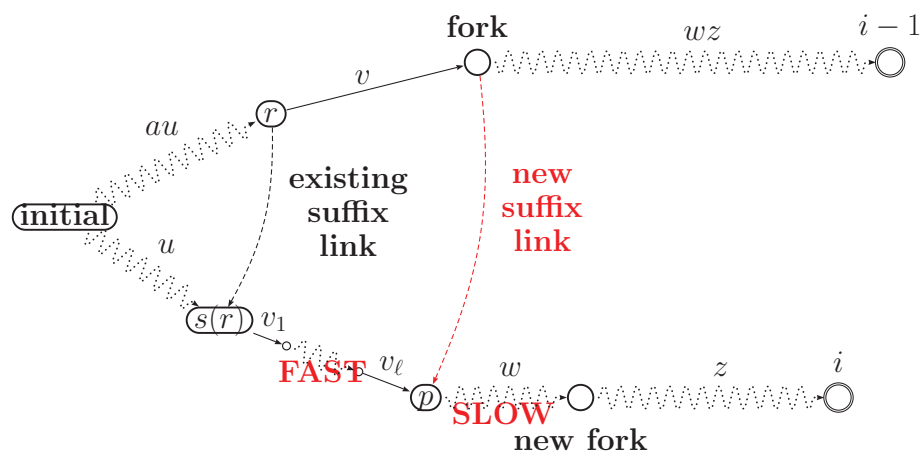
- ★ Scheme for the insertion of suffix $y[i \dots n-1] = u \cdot v \cdot w \cdot z$



- ★ It first computes $p = \text{TARGET}(s[r], v)$ with FAST-FIND (if necessary)
- ★ then the fork of the current suffix with SLOW-FIND

Scheme for insertion (continued)

- ★ General scheme for inserting the next suffix in the data structure when the suffix target of the current fork is not defined



Complete algorithm

SUFFIX-TREE(y)

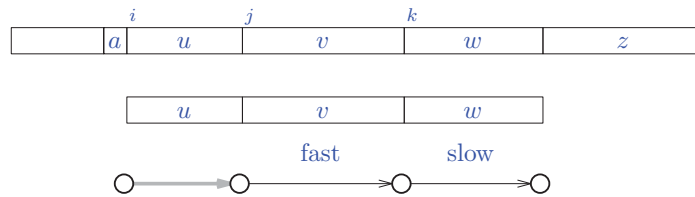
```

1   $T \leftarrow \text{NEW-TREE}()$ 
2   $s[\text{initial}[T]] \leftarrow \text{initial}[T]$ 
3   $(\text{fork}, k) \leftarrow (\text{initial}[T], 0)$ 
4  for  $i \leftarrow 0$  to  $n - 1$  do
5       $k \leftarrow \max\{k, i\}$ 
6      if  $s[\text{fork}] = \text{NIL}$  then
7           $r \leftarrow \text{parent of fork}$ 
8           $(j, \ell) \leftarrow \text{label}(r, \text{fork})$ 
9          if  $r = \text{initial}[T]$  then
10              $\ell \leftarrow \ell - 1$ 
11              $s[\text{fork}] \leftarrow \text{FAST-FIND}(s[r], k - \ell, k)$ 
12              $(\text{fork}, k) \leftarrow \text{SLOW-FIND}(s[\text{fork}], k)$ 
13             if  $k < n$  then
14                  $q \leftarrow \text{NEW-STATE}()$ 
15                  $\text{Adj}[\text{fork}] \leftarrow \text{Adj}[\text{fork}] \cup \{(k, n - k), q\}$ 
16                  $\text{output}[q] \leftarrow i$ 
17             else  $\text{output}[\text{fork}] \leftarrow i$ 
18  $\text{output}[\text{initial}] \leftarrow n$ 
19 return  $T$ 

```

Running time

★ Scheme for insertion



- ★ Main iteration increments i , which never decreases
- ★ Iteration in FAST-FIND increments j , which never decreases
- ★ Iteration in SLOW-FIND increments k , which never decreases
- ★ Basic operations run in constant time or in time $O(\log \text{card } \Sigma)$

Theorem 1 *Execution of SUFFIX-TREE(y) = $\mathcal{S}(y)$ takes $O(|y| \times \log \text{card } \Sigma)$ time in the comparison model.*