

# 6. Linear-Time Construction of Suffix Trees McCreight's Algorithm

Kena Alexander 26/05/2015

# McCreight's Method

---

- **McCreight's method**
  - Linear-time construction algorithm
  - Space-saving improvement over Weiner's method.
  - Same space efficiency as Ukkonen's method

# Suffix Trees

- **Definition**

- A suffix tree of a sequence,  $x$  is a *compressed trie* of all suffixes of the sequence  $x\$$
- $\$$  is special end character

$x = a\ b\ a\ a\ b\ \$$

1     $a\ b\ a\ a\ b\ \$$

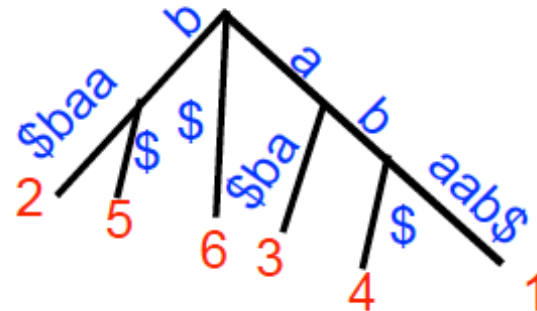
2        $b\ a\ a\ b\ \$$

3            $a\ a\ b\ \$$

4                $a\ b\ \$$

5                    $b\ \$$

6                        $\$$



# McCreight's Algorithm at a High Level

---

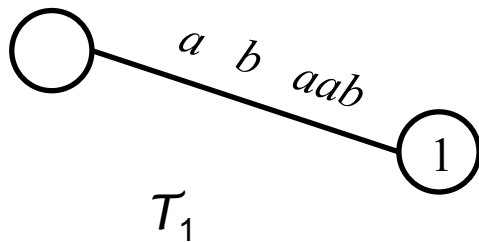
- **McCreight's algorithm**
  - Constructs the suffix tree  $T$  for  $m$ -length string  $S$  by
    - Inserting suffixes in order, one at a time starting from suffix one.
    - Iteratively for  $i=1, \dots, m+1$  build tries,  $T_i$ ,
    - where each in-between iteration is a  $T_i$  trie of sequences  $x[1..m+1], x[2..m+1], \dots, x[i..m+1]$

# McCreight's Algorithm

$i$	1	2	3	4	5	6
$S$	$a$	$b$	$a$	$a$	$b$	$\$$

- **Example**

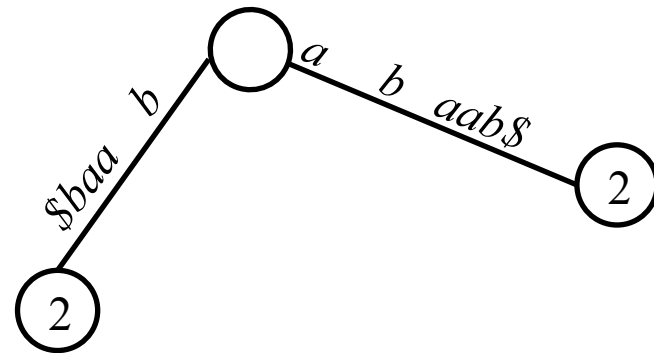
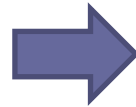
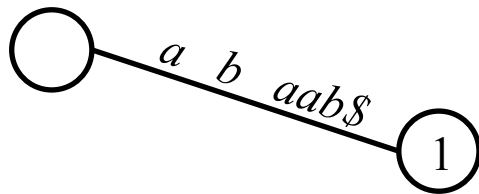
- Divided into  $m$  phases
  - In phase  $i + 1$ , tree  $\mathcal{T}_{i+1}$  is constructed from  $\mathcal{T}_i$



# McCreight's Algorithm

- Example**

$i$	1	2	3	4	5	6
$S$	$a$	$b$	$a$	$a$	$b$	$\$$

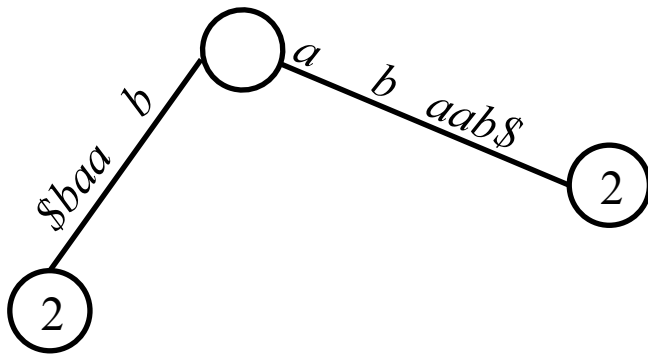


$T_2$

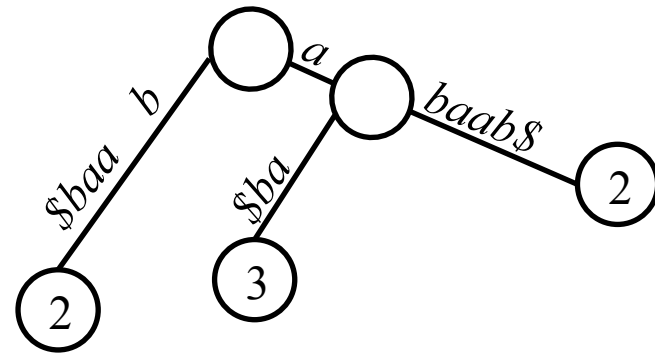
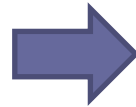
# McCreight's Algorithm

- Example**

$i$	1	2	3	4	5	6
$S$	$a$	$b$	$a$	$a$	$b$	$\$$



$T_2$

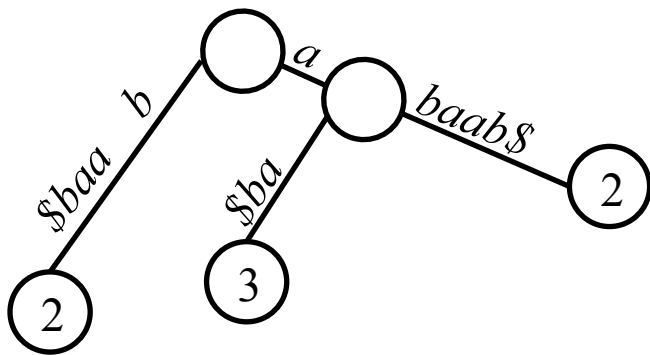


$T_3$

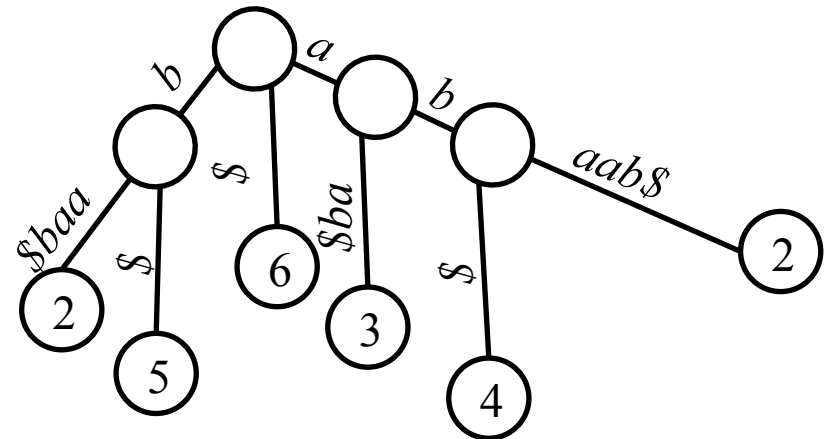
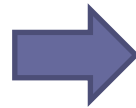
# McCreight's Algorithm

- Example**

$i$	1	2	3	4	5	6
$S$	$a$	$b$	$a$	$a$	$b$	$\$$



$T_3$



$T_{m+1}$



# McCreight's Algorithm at a High Level

---

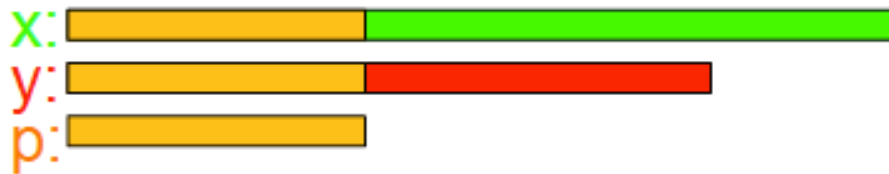
- **McCreight's algorithm**
  - Essential Trick
    - The essential trick is being clever in how we insert  $x[i..m]$  into  $T_i$  so we don't spend  $O(m^2)$  time to build the tree.

# Terminology

---

- **Common Prefix**

- A common prefix of strings x and y is a string p that is a prefix of both.



- **Example:**

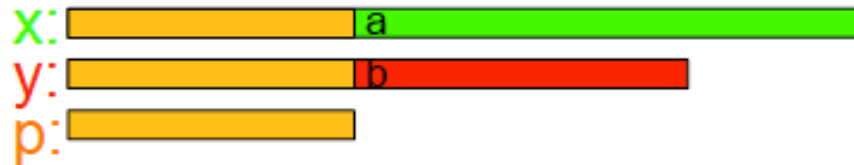
- $x = \text{abcabxabcd}$
- $y = \text{abcdabxa}$
- $p = \text{abc}$

Note: a, ab and abc are all common prefixes.

# Terminology

- **Longest Common Prefix**

- The longest common prefix,  $p = \text{LCP}(x, y)$ , is a prefix such that:  $x[|p|+1] \neq y[|p|+1]$

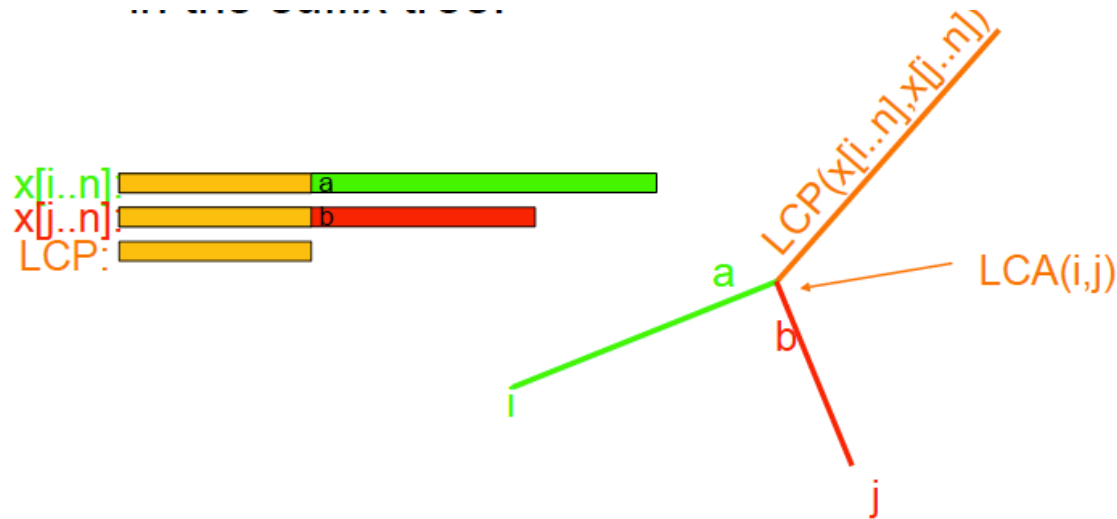


- **Example**

- $x = \text{abcabxabc}$
- $y = \text{abcdabxa}$
- $\text{lcp} = \text{abc}$

# Terminology

- **Lowest Common Ancestor**
- For suffixes of  $x$ ,  $x[i..n]$ ,  $x[j..n]$ , their longest common prefix is their lowest common ancestor in the suffix tree:



# Head and Tail

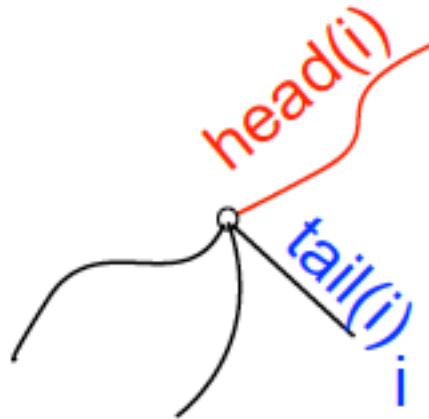
---

- **Head**

- Let  $\text{head}(i)$  denote the longest LCP of  $x[i..n]$  and  $x[j..n]$  for all  $j < i$

- **Tail**

- Let  $\text{tail}(i)$  be the string such that  $x[i..n] = \text{head}(i)\text{tail}(i)$



# Head and Tail

---

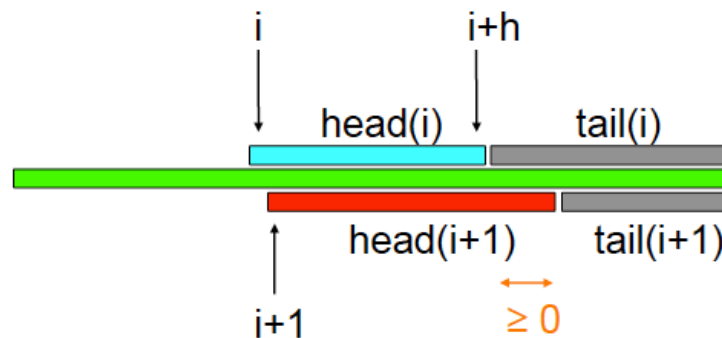
- **Iteration**

- Iteration  $i$  in McCreight's algorithm consist of
  - finding (or inserting) the node for head( $i$ ),
  - and appending tail( $i$ ))
- The trick is a clever way of finding head( $i$ )

# Lemma

- **Lemma**


- Let  $\text{head}(i) = x[i..i+h]$ . Then  $x[i+1..i+h]$  is a prefix of  $\text{head}(i+1)$



- $x = \text{abcabxabcd}$
- $\text{Head} = x[1..3] = \text{abc}$
- $\text{Head}+1 = \text{bcd}$
- $x[2..3] = \text{bc}$

# Lemma

- **Proof**

- Trivial for  $h=0$  ( $\text{head}(i)$  empty), so assume  $h>0$ :
- Let  $\text{head}(i) = ay$ :  $|a$  
- By def. exists  $j<i$  such that  $\text{LCP}(i,j)=ay$
- Thus suffix  $j+1$  and  $i+1$  share prefix  $y$
- Thus  $y$  is a prefix of  $\text{LCP}(i+1,j+1)$
- Thus  $y$  is a prefix of  $\text{head}(i+1)$

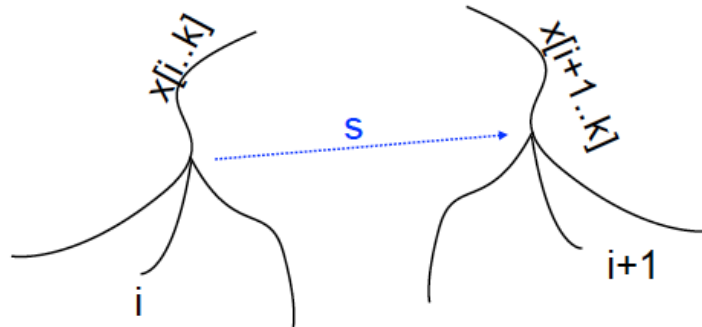




# Suffix link

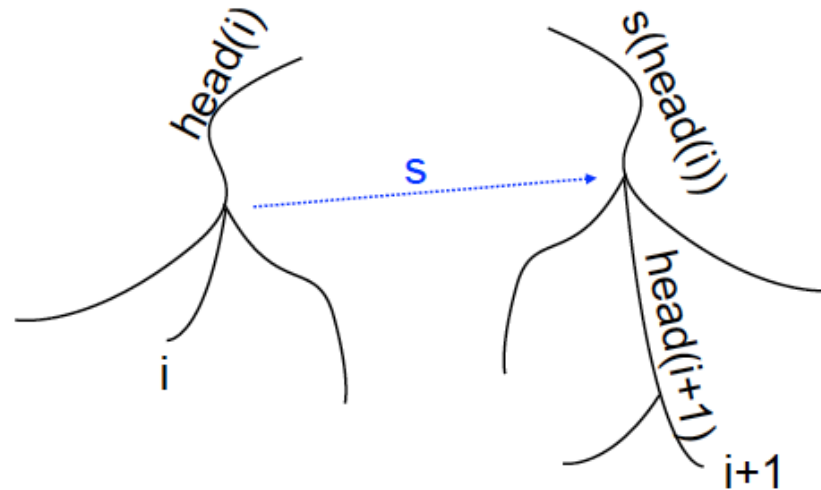
- **Definition of *suffix link***

- Let  $u = x\alpha$  denote an arbitrary string, where  $x$  denotes a single character and  $\alpha$  denotes a (possibly empty) substring.
- Then the suffix link  $s(u) = \alpha$  for  $u = x\alpha$  or
- $s(u) = ""$  (the root node) if  $u = ""$
- Suffix link is a pointer from a path labeled  $x[i..k]$  to a path labeled  $x[i+1..k]$ :



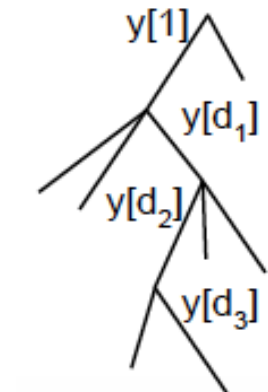
# Corollary of Lemma

- $s(\text{head}(i))$  is a prefix of  $\text{head}(i+1)$
- Thus  $s(\text{head}(i))$  is an ancestor of  $\text{head}(i+1)$



# Slowscan and Fastscan

- Slowscan:
  - if we do not know if string  $y$  is in  $T_i$ , we must search character by character
- Fastscan (Rescanning):
  - if we do know that  $y$  is in  $x$ , we can jump directly from node to node
    - At node  $u$  at (path-depth)  $d$ , follow the edge with label starting with  $y[d]$
    - Continue until we reach the end of  $y$ .
    - We will either end on a
      - Node (if  $y$  is in  $T_i$ )
      - Edge (if  $y$  is a prefix of a string in  $T_i$ )
        - $s(\text{head}(i))$  is a prefix of  $\text{head}(i+1)$
        - Thus  $s(\text{head}(i))$  is an ancestor of  $\text{head}(i+1)$



# Sketch of McCreight's Algorithm

---

- Begin with the tree  $T_1$ :
- For  $i=1, \dots, n$ , build tree  $T_{i+1}$  satisfying:
  - $T_{i+1}$  is a compressed trie for  $x[j..n]$ ,  $j \leq i+1$
  - All non-terminal nodes (with the possible exception of  $\text{head}(i)$ ) have a suffix link  $s(-)$
- Each iteration must:
  - Add node  $i+1$
  - Potentially add  $\text{head}(i+1)$
  - Add  $\text{tail}(i+1)$
  - Add suffix link  $\text{head}(i) \rightarrow s(\text{head}(i))$

# iteration $i$

---

- Beginning of iteration
  - Let  $\text{head}(i) = uv$
  - $\text{parent}(\text{head}(i)) = u$
  - $w = s(u)v = s(\text{head}(i))$
- Rest of iteration
  - Move quickly to  $w$ , then search for  $\text{head}(i+1)$  starting there
    - By the invariant,  $s(\text{parent}(\text{head}(i)))$  and the suffix link exists;
    - by the lemma,  $w$  is an ancestor of  $\text{head}(i+1)$

# Observe: $w$ is in $T$

---

- $\text{head}(i)$  is a prefix of  $x[j..n]$  for some  $j < I$
- Thus  $w$  is a prefix of  $x[j+1..n]$  for some  $j < I$ 
  - i.e  $w$  is a prefix of some suffix  $j \leq i$
  - i.e.  $w$  is in  $T_i$
- Consequently: we can search for  $w$  from  $s(u)$  using fastscan!

# W

---

- If  $w$  is a node
  - Update  $s(\text{head}(i)) := w$
  - Then search for  $\text{head}(i+1)$  using slowscan
- If  $w$  is on an edge
  - If  $w$  is not a node, then all suffix  $j < i$  with prefix  $w$  agree on the next letter
  - By definition of  $\text{head}(i)$  there is  $j < i$  such that suffix  $x[i..n]$  and  $x[j..n]$  differs after  $\text{head}(i)$ 
    - $x[i+1..n]$  must also disagree at that character
    - Thus  $\text{head}(i+1)$  must be  $w$
- Add node  $w$ , update  $\text{head}(i) := w$  and set  $\text{head}(i+1) = w$

# McCreight's Algorithm

---

- Insert suffixes  $x[1..m] \rightarrow x[m..m]$
- Step 1
  - Walk to the deepest node with a suffix link
- Step 2
  - Traverse down  $\text{head}(i+1)$  by looking at the first character of the edge
- Step 3
  - Traverse the rest of the suffix until we encounter a mismatch
  - Create a new node and using the remaining characters add an edge.



# McCreight's Step 1

---

- Find the deepest node with a suffix link (lets call this head(i))
- head(i) can be divided into three substrings
  - $x\alpha\beta$
  - $x\alpha$  is the path label of the node y. if no such node exists  $x\alpha = ""$
  - If such a node exists then  $x$  is the first character and  $\alpha$  is the rest of the characters in the string.
- If we follow the suffix link  $s(\mathbf{x\alpha})$  we get to w whose label is  $\alpha$
- The rest of the string will be  $\beta$

# McCreight's Step 2

---

- Rescanning phase {Fastscan}
- Rescan the characters of  $\beta$  from node  $w$ 
  - (analyze only the first characters and jump from node to node until a mismatch or no nodes to jump)
- We know that  $\beta$  is part of the tree therefore
  - The number of steps required to rescan is proportional to the intermediate nodes visited while traversing  $\beta$
- If there is no node at this point, create a new node,  $d$
- Update the suffix link  $s(\text{head}(i))$  to point to  $d$

# McCreight's Step 3

---

- Scanning phase {Slowscan}
- Let the remaining part of  $\text{head}(i) = \text{tail}(i)$ . since we do not know the length of this we have to analyze each character down the subtree.
- If we get a mismatch then
  - Create a new node and let the remainder of characters be the leaf from this node.

# Example

<i>i</i>	1	2	3	4	5	6
<i>S</i>	a	b	a	a	b	\$

**Algorithm** McCreight;

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}(i)$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

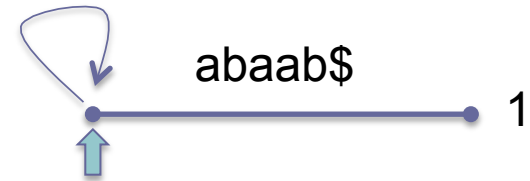
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

<i>i</i>	1	2	3	4	5	6
<i>S</i>	a	b	a	a	b	\$

**Algorithm** McCreight;

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}()$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$

jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$

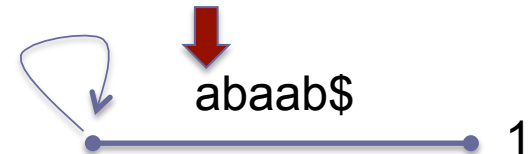
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

<i>i</i>	1	2	3	4	5	6
<i>S</i>	a	b	a	a	b	\$

**Algorithm** McCreight;

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}(i)$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

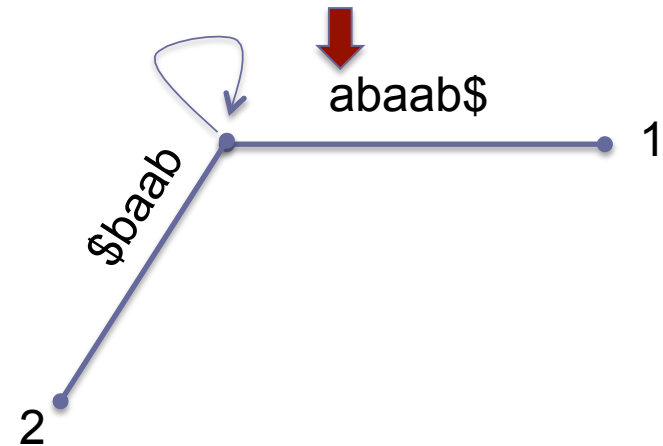
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm** McCreight;

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}(i)$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

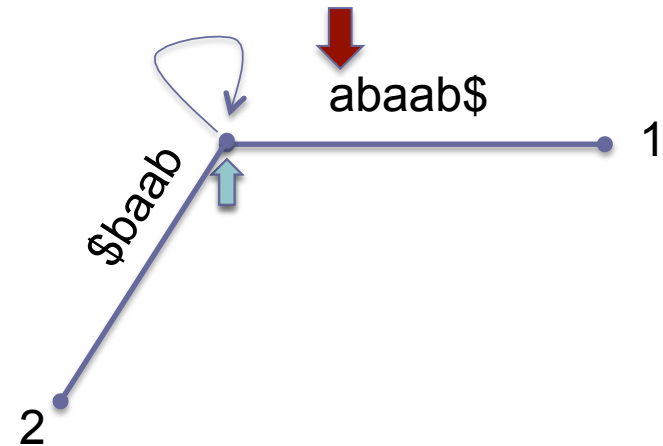
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm McCreight;**

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}(i)$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$

jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

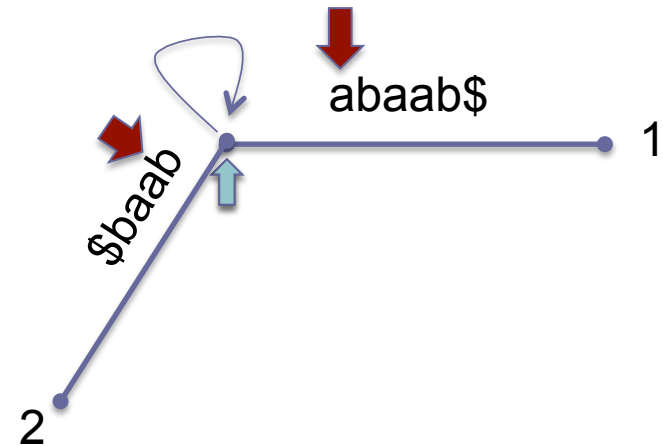
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.





# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm** McCreight;

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}()$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

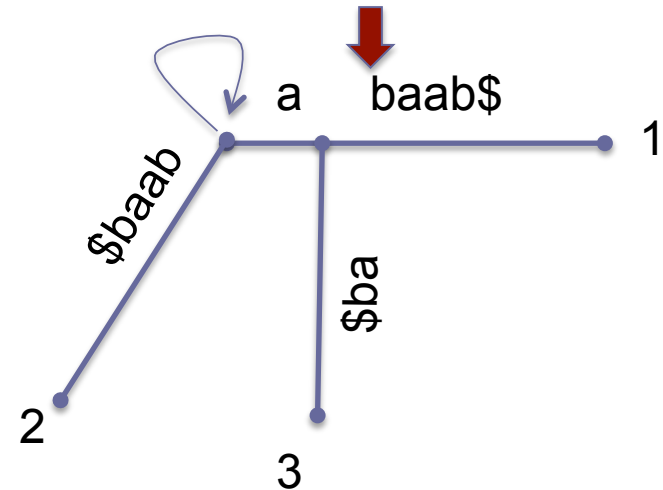
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm McCreight;**

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}(i)$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

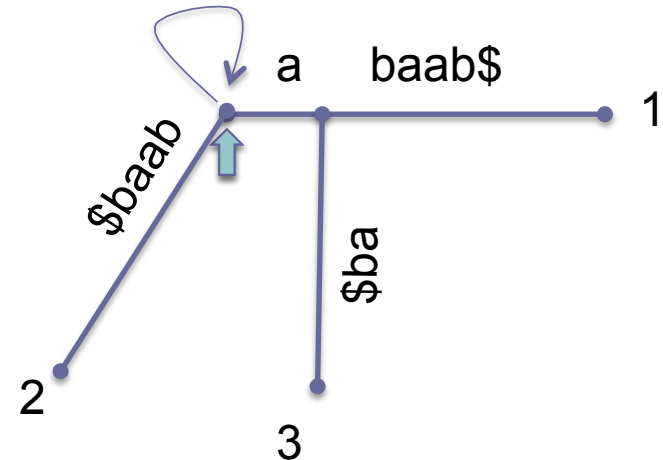
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm McCreight;**

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}()$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$

jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

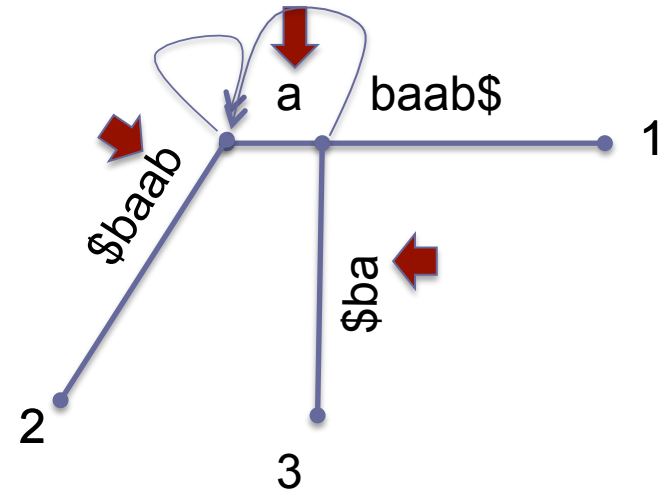
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm** McCreight;

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}()$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

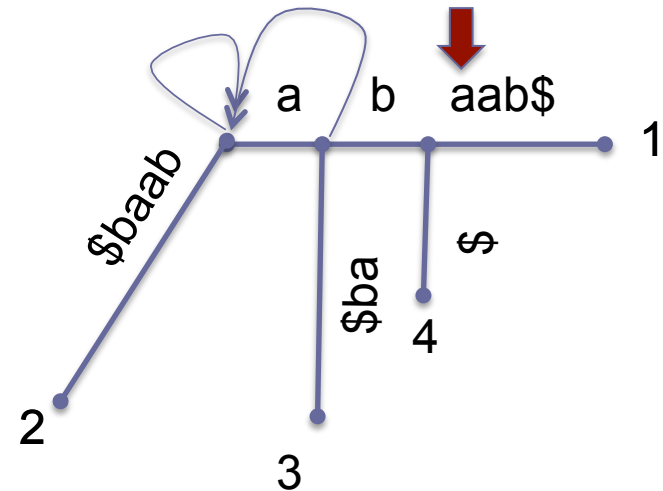
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm** McCreight;

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}(i)$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

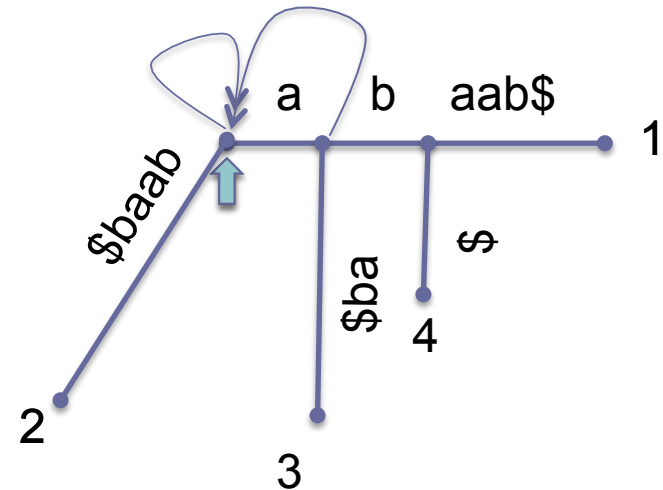
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

**Algorithm McCreight;**

**Step1:**

Find the deepest node with a suffix link

Follow the suffix link  $\text{suff}(i)$  to get to  $\text{head}(i) = \alpha\beta$

**Step2: {Fastscan}**

analyze the first characters of  $\text{head}(i)$  comparing to  $\beta$   
jumping nodes if match.

if there is no node then

create a node =  $d$

update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$

**Step3: {Slowscan}**

let the remainder of  $\text{head}(i)$  be  $\text{tail}(i)$

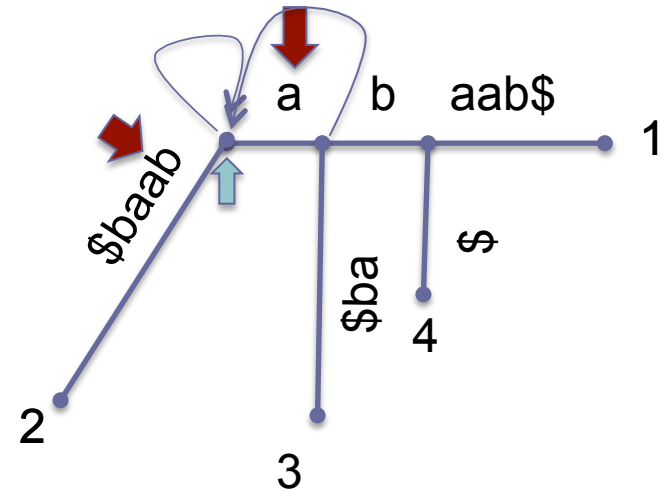
compare each character of  $\text{tail}(i)$  to rest of characters of  $\beta$   
until a mismatch.

if we get a mismatch

create a new node.

let the remaining characters be the leaf.

end.



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$

### Algorithm McCreight;

### Step1:

Find the deepest node with a suffix link

**Follow the suffix link  $\text{suff}()$  to get to  $\text{head}(i) = \alpha\beta$**

## Step2: {Fastscan}

**analyze the first characters of head(i) comparing to  $\beta$   
jumping nodes if match.**

if there is no node then

**create a node = d**

**update the suffix link  $\text{suff}(\text{head}(i)) \rightarrow d$**

### Step3: {Slowscan}

let the remainder of head(i) be tail(i)

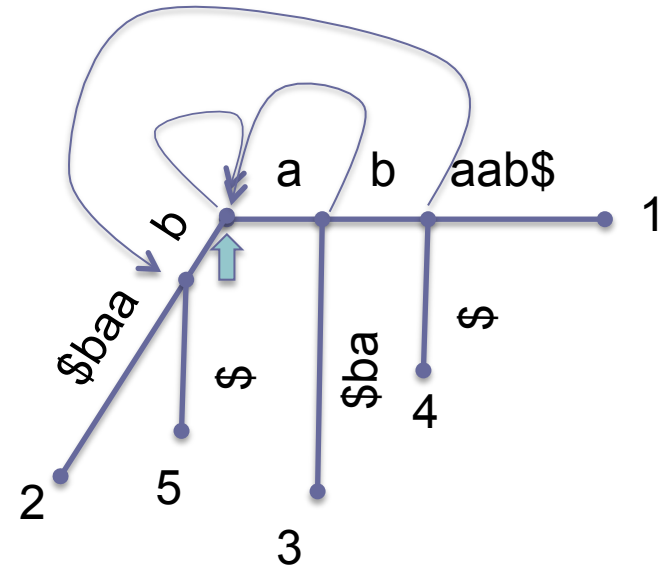
compare each character of tail(i) to rest of characters of  $\beta$   
until a mismatch.

**if we get a mismatch**

**create a new node.**

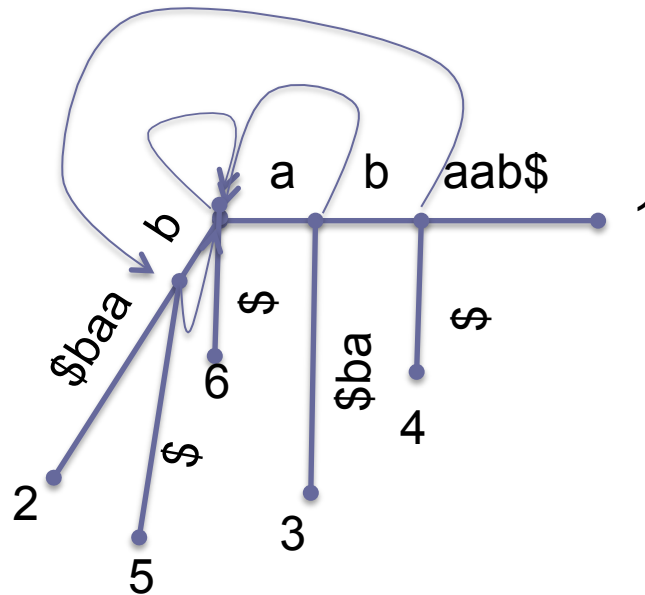
**let the remaining characters be the leaf.**

**end.**



# Example

$i$	1	2	3	4	5	6
$S$	a	b	a	a	b	\$





# Time analysis

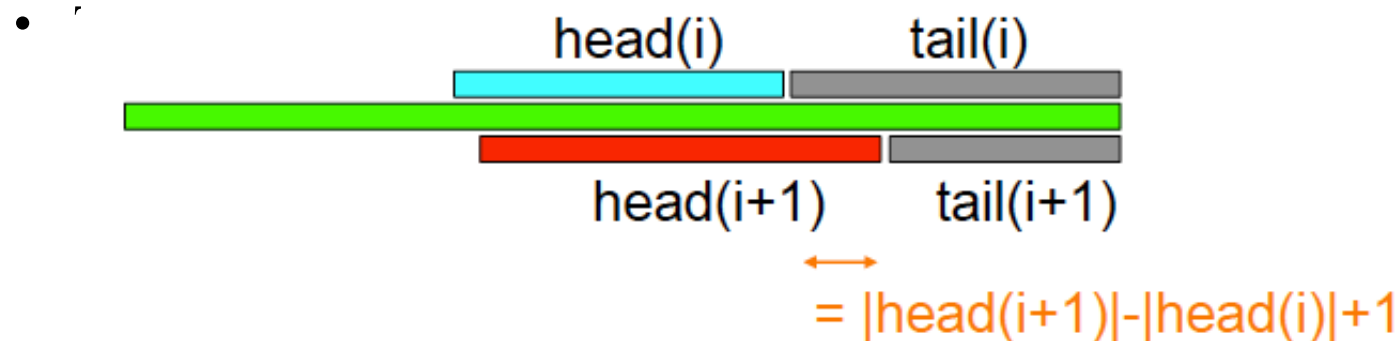
---

- Everything but searching is done in constant time per suffix,
- The running time is  $O(n + \text{“slowscan”} + \text{“fastscan”})$ .

# Time analysis

- Slowscan

- We use slowscan to find  $\text{head}(i+1)$  from  $w=s(\text{head}(i))$ ,
- The complexity of one run of *slowscan* at stage  $i$  is proportional to  $|\text{head}(i+1)| - |\text{head}(i)| + 1$



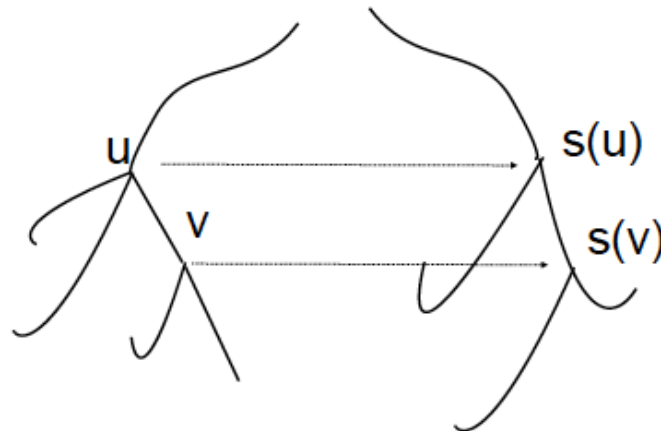
# Time analysis

---

- Fastscan
  - Fastscan uses time proportional to the number of nodes it processes.
  - *If we define  $d(v)$  as the depth of node  $v$* 
    - Fastscan increases the node depth
    - Following parent and suffix pointers decreases the node depth
- Time usage of fastscan is bounded by the total depth-increase (amortized analysis)

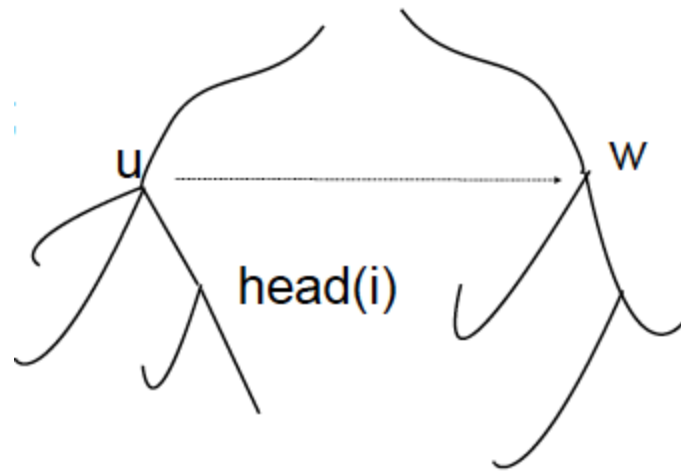
# Time analysis

- Proposition
  - $d(v) \leq d(s(v)) + 1$
- Proof:
  - For any ancestor  $u$  of  $v$ ,  $s(u)$  is an ancestor of  $s(v)$
  - Except for the empty prefix and the single letter prefix of  $v$ , the  $s(u)$ 's are different



# Time analysis

- Corollary
  - In each step, before calling fastscan, we decrease the depth by at most 2:
    - $d(u) = d(\text{head}(i)) - 1$ ;
    - $d(w) \geq d(u) - 1$
  - The total decrease is thus  $2n$



# Time analysis

---

- The time usage of fastscan is bounded by  $n$  plus the total decrease of depth,
  - i.e. the time usage of fastscan is  $O(n)$

# Summary

---

- We iteratively build tries of suffixes of  $x$ .
- Using suffix links and fastscan we can quickly find where to insert the next suffix in our current trie.
- By amortized analysis, the total running time becomes linear.