

# On-Line Construction

## of Suffix Trees

E. Ukkonen

1

## Overview

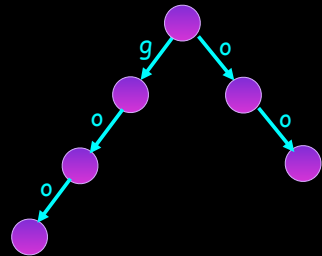
- Suffix tries
- On-line construction of suffix tries in quadratic time
- Suffix trees
- On-line construction of suffix trees in linear time
- Applications

2

## Suffix Trees

A **suffix tree** is a trie-like data structure representing all suffixes of a string.

goo



3

## Notations

- Let  $T = t_1 \dots t_n$  be a string.
- For  $0 \leq i \leq n$ , let  $T^i = t_1 \dots t_i$  denote the  $i$ -length prefix of  $T$ .
- For  $1 \leq i \leq n+1$ , let  $T_i = t_i \dots t_n$  denote the suffix of  $T$  that starts at the  $i^{\text{th}}$  position.
- Let  $\sigma(T) = \{T_i \mid 1 \leq i \leq n+1\}$ .

4



## The Size of Suffix Tries

**Theorem:** The size of  $STrie(T)$ , where  $|T| = n$ , is  $O(n^2)$ .

**Proof:** The size of  $STrie(T)$  is linear in the number of substrings of  $T$ .  
 $T$  has at most  $O(n^2)$  substrings. Thus the size of  $STrie(T)$  is  $O(n^2)$ .

9

## On-Line Construction of Suffix Tries

- Let  $T = t_1 \dots t_n$ .
- $\forall 1 \leq i \leq n$ , the algorithm constructs  $STrie(T^i)$ .
- First we construct  $STrie(T^0) = STrie(\varepsilon)$ .
- Then,  $\forall 1 \leq i \leq n$ , we obtain  $STrie(T^i)$  from  $STrie(T^{i-1})$ .

10

## On-Line Construction of Suffix Tries (cont.)

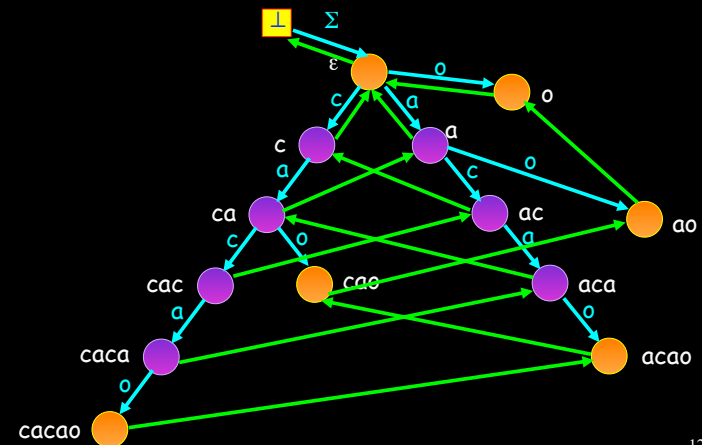
**Observation 1:**  $\sigma(T^i) = \{xt_i \mid x \in \sigma(T^{i-1})\} \cup \{\varepsilon\}$ .

**Observation 2:** The suffixes of  $T^i$  can be found by starting at the state  $T^i$  and following the suffix links, until  $\varepsilon$ .  
 Thus,  $\sigma(T^i) = \{f^j(T^i) \mid 0 \leq j \leq i\}$ .

**Definition:** The path from  $T^i$  to  $\perp$  following the suffix links is called the **boundary path** of  $STrie(T^i)$ .

11

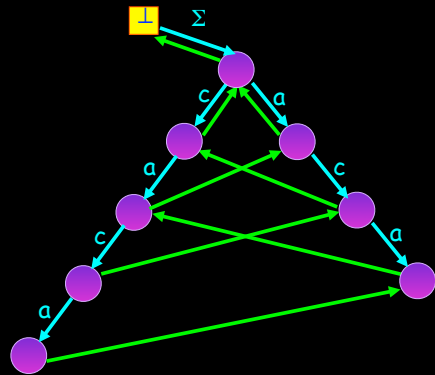
## On-Line Construction of Suffix Tries (cont.)



12

## STrie( $T^{i-1}$ ) $\Rightarrow$ STrie( $T^i$ )

cac  $\Rightarrow$  caca



13

## The Algorithm

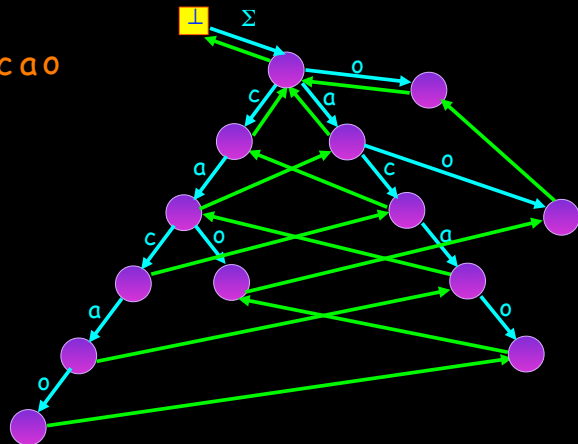
```

create STrie( $\epsilon$ )
top  $\leftarrow \epsilon$ 
for  $i \leftarrow 1$  to  $n$  do
   $r \leftarrow \text{top}$ 
  while  $g(r, t_i)$  is undefined do
    create new state  $r'$  and  $g(r, t_i) \leftarrow r'$ 
    if  $r \neq \text{top}$  then  $f(\text{old-}r') \leftarrow r'$ 
     $\text{old-}r' \leftarrow r'$ 
   $r \leftarrow f(r)$ 
 $f(\text{old-}r') \leftarrow g(r, t_i)$ 
top  $\leftarrow g(\text{top}, t_i)$ 
    
```

14

## The Algorithm (cont.)

cacao



15

## Running Time

**Theorem:** The running time of the algorithm is linear in the size of  $\text{STrie}(T)$ , which is, in worst case,  $O(|T|^2)$ .

16

## Running Time (cont.)

```
create STrie( $\epsilon$ )
top  $\leftarrow \epsilon$ 
for  $i \leftarrow 1$  to  $n$  do
   $r \leftarrow$  top
  while  $g(r, t_i)$  is undefined do
    create new state  $r'$  and  $g(r, t_i) \leftarrow r'$ 
    if  $r \neq$  top then  $f(\text{old-}r') \leftarrow r$ 
    old- $r' \leftarrow r'$ 
     $r \leftarrow f(r)$ 
   $f(\text{old-}r') \leftarrow g(r, t_i)$ 
```

$O(1)$  for each  
node added to  
STrie(T)

17

## Suffix Trees

- A **suffix tree**  $STree(T)$  represents  $STrie(T)$  in space linear in  $|T|$ .
- This is achieved by representing only a subset of  $Q \cup \{\perp\}$  of  $QU\{\perp\}$ , called the **explicit states**.

18

## Explicit and Implicit States

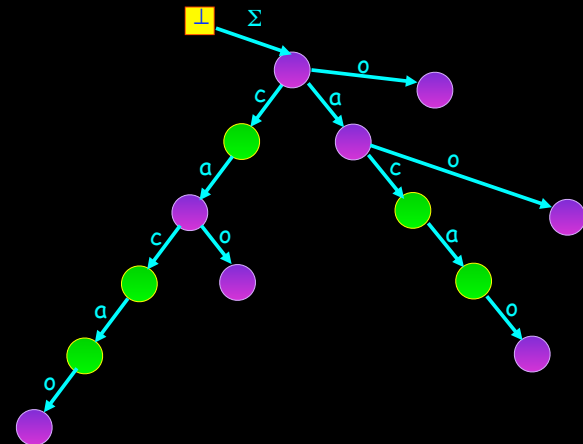
**Definition:** A state  $q$  is called **explicit** in the following cases:

- $q$  is a leaf
- $q$  is a branching state (has at least two transitions)
- **root** and  $\perp$  are also defined to be branching states.

Otherwise (if  $q$  has exactly one transitions and is not the **root** or  $\perp$ ),  $q$  is called **implicit**.

19

## Explicit and Implicit States (cont.)

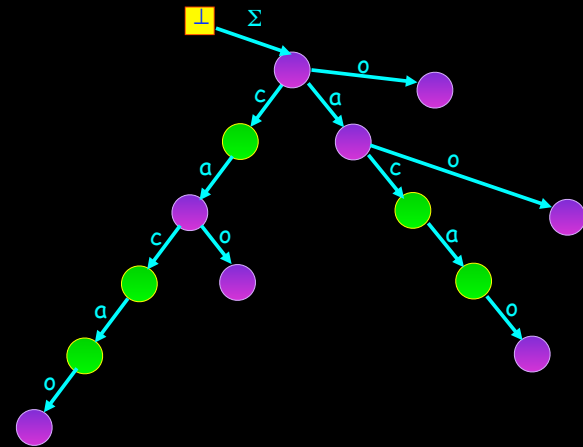


20

## Generalized Transition Function

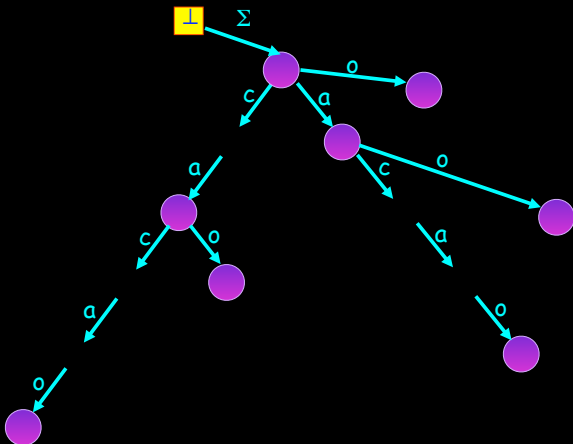
- The string  $w$  spelled out by the transition path in  $STrie(T)$  between two explicit states  $s$  and  $r$  is represented in  $STree(T)$  as a **generalized transition**  $g'(s,w) = r$ .
- A generalized transition  $g'(s,w) = r$  is called an **a-transition** if  $\exists a \in \Sigma$  and  $v \in \Sigma^*$  s.t.  $w = av$ .
- Note that for each explicit state  $s$  and  $a \in \Sigma$  there is at most one **a-transition** from  $s$ .

## $STrie(T) \Rightarrow STree(T)$



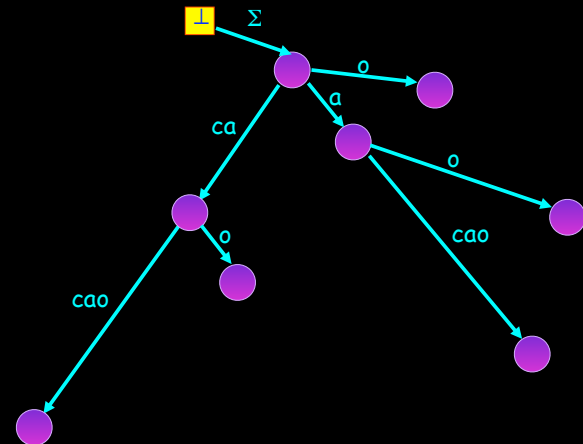
22

## $STrie(T) \Rightarrow STree(T)$



23

## $STrie(T) \Rightarrow STree(T)$



24

## Suffix Links

**Definition:** If  $x \in Q'$  is a branching state and  $x = ay$ , where  $a \in \Sigma$ , then the **suffix link** of  $x$  is defined by  $f'(x) = y$ , and  $f'(\epsilon) = \perp$ .

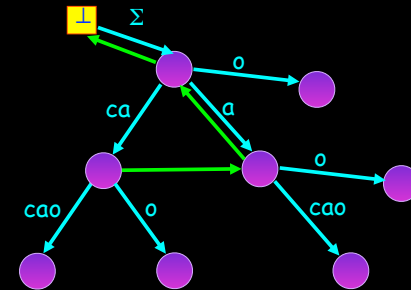
**Proposition:** If  $x \in Q'$  is a branching state and  $f'(x) = y$  then  $y$  is also a branching state.

**Proof:**  $\exists a \neq b \in \Sigma$  s.t.  $xa$  and  $xb$  are substrings of  $T$ .  $y$  is a suffix of  $x$ . Thus  $ya$  and  $yb$  are also substrings of  $T$ .

25

## STree(T)

$S\text{Tree}(T) = (Q' \cup \{\perp\}, \text{root}, g', f')$ .



26

## The Size of Suffix Trees

**Theorem:** The size of  $S\text{Tree}(T)$ , where  $|T| = n$ , is  $O(n)$ .

**Proof:** Since we represent each substring  $w = t_k \dots t_p$  of  $T$  by a pair pointers  $(k, p)$ , the size of  $S\text{Tree}(T)$  is linear in the number of explicit states.  $S\text{Tree}(T)$  has at most  $n$  leaves, and thus at most  $n - 1$  branching states. Therefore, the size of  $S\text{Tree}(T)$  is  $O(n)$ .

27

## Reference Pairs

**Definition:** Let  $r$  be an explicit or implicit state.  $(s, w)$  is called a **reference pair** for  $r$  if:

- $s$  is an explicit state and an ancestor of  $r$ .
- $w$  is the string spelled out by the transitions from  $s$  to  $r$  in the corresponding suffix trie.

**Definition:** A reference pair  $(s, w)$  for  $r$  is called **canonical** if  $s$  is the closest explicit ancestor of  $r$  (or  $r$  itself, if it is explicit).

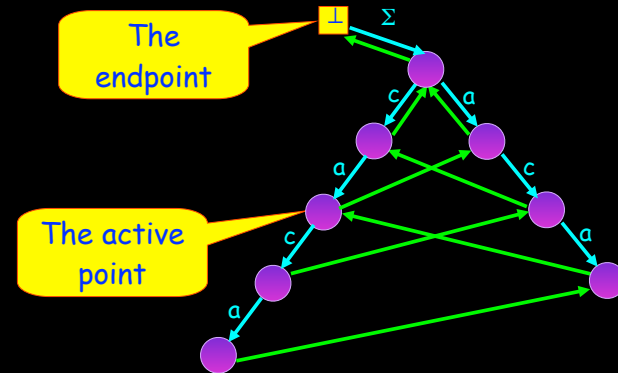
## Active Point and Endpoint

Let  $s_1 = T^{i-1}$ ,  $s_2, \dots, s_i = \text{root}$ ,  $s_{i+1} = \perp$  be the boundary path of  $\text{STrie}(T^{i-1})$ .

**Definition:**  $s_j$  is called the **active point** of  $\text{STrie}(T^{i-1})$  if  $j$  is the smallest index for which  $s_j$  is not a leaf.

**Definition:**  $s_{j'}$  is called the **endpoint** of  $\text{STrie}(T^{i-1})$  if  $j'$  is the smallest index for

## Active Point and Endpoint (cont.)



30

## Active Point and Endpoint (cont.)

**Proposition:**  $s_j$  and  $s_{j'}$  are well defined and  $j \leq j'$ .

**Proof:**

- $\text{root}$  is not a leaf  $\Rightarrow s_j$  is defined.
- $g(\perp, t_i)$  is defined  $\Rightarrow s_{j'}$  is defined.
- $g(s_j, t_i)$  is defined  $\Rightarrow s_{j'}$  is not a leaf  $\Rightarrow j \leq j'$ .

31

## Adding $t_i$ -Transitions to $\text{STrie}(T^{i-1})$

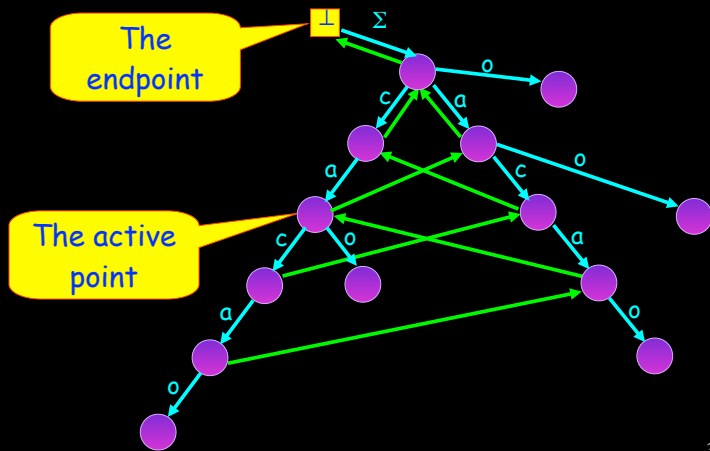
**Lemma:** When obtaining  $\text{STrie}(T^i)$  from  $\text{STrie}(T^{i-1})$  the algorithm adds a  $t_i$ -transition to each state  $s_h$  s.t.  $1 \leq h < j'$ , and only to these states, as follows:

- For  $1 \leq h < j$ , the new transition expands an old branch of the trie that ends at  $s_h$ .
- For  $j \leq h < j'$ , the new transition initiates a new branch from  $s_h$ .

32



## Adding $t_i$ -Transitions to $STrie(T^{i-1})$ (cont.)



33

## On-Line Construction of Suffix Trees

- We create  $STree(\varepsilon)$ , and then  $\forall 1 \leq i \leq n$  we obtain  $STree(T^i)$  from  $STree(T^{i-1})$ .
- When obtaining  $STree(T^i)$  from  $STree(T^{i-1})$ , we update  $STree(T^{i-1})$  according to the transitions we would add to  $STrie(T^{i-1})$ .
- Note that  $s_1, \dots, s_{i-1}$  are not necessarily explicit states.

34

## On-Line Construction of Suffix Trees (cont.)

For  $1 \leq h < j$ :

- $s_h$  is a leaf. Thus,  $\exists s, 0 \leq k \leq i-1$  s.t.  $g'(s, (k, i-1)) = s_h$ . We replace this transition by  $g'(s, (k, i)) = s_h$ .
- This would take too much time. Thus, we denote transitions of the type  $g'(s, (k, i-1))$  in  $STree(T^{i-1})$  by  $g'(s, (k, \infty))$ . Hence, no updates are needed.

35

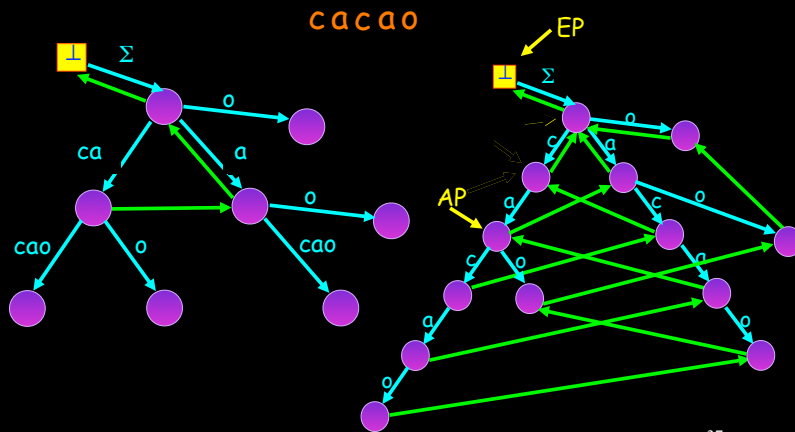
## On-Line Construction of Suffix Trees (cont.)

For  $j \leq h < j'$ :

- If  $s_h$  is an implicit state, we turn it into an explicit state by splitting the transition containing it.
- We create a new leaf  $s_h t_i$  and add a new transition  $g'(s_h, (i, \infty))$ .

36

## On-Line Construction of Suffix Trees (cont.)



37

## Lemma 1

Lemma 1: Let  $(s, (k, p))$  be some reference pair for a state  $r$ . Then  $\exists s', k'$  s.t.  $(s', (k', p))$  is the canonical reference pair for  $r$ .

Proof: Let  $s'$  be the closest explicit ancestor of  $r$ , or  $r$  itself if  $r$  is explicit.  $t_k \dots t_p$  is the path from the explicit state  $s$  to  $r$ . Thus, the path from  $s'$  to  $r$  is a suffix  $t_{k'} \dots t_p$  of  $t_k \dots t_p$ .

38

## Lemma 2

Lemma 2: Let  $r$  be a state on the boundary path of  $STrie(T^i)$ . Then  $\exists s, k$  s.t.  $(s, (k, i))$  is the canonical reference pair for  $r$ .

Proof:  $r$  is on the boundary path of  $STrie(T^i)$ .

- $\Rightarrow r$  refers to some suffix  $t_k \dots t_i$  of  $T^i$ .
- $\Rightarrow (\varepsilon, (k, i))$  is a reference pair for  $r$ .
- $\Rightarrow$  the claim holds by lemma 1.

39

## Lemma 3

Lemma 3: Let  $(s, (k, i-1))$  be a reference pair for the endpoint of  $STrie(T^{i-1})$ . Then  $(s, (k, i))$  is a reference pair for the active point of  $STrie(T^i)$ .

Proof:

- $s_j$  is the active point of  $STrie(T^{i-1})$  iff  $t_j \dots t_{i-1}$  is the longest suffix of  $T^{i-1}$  that occurs at least twice in  $T^{i-1}$ .

40

## Lemma 3 (cont.)

### Proof (cont.):

- $s_j$  is the endpoint of  $\text{STrie}(T^{i-1})$  iff  $t_j \dots t_{i-1}$  is the longest suffix of  $T^{i-1}$  such that  $t_j \dots t_{i-1} t_i$  is a substring of  $T^{i-1}$ .
- Thus, if  $s_j$  is the endpoint of  $\text{STrie}(T^{i-1})$ , then  $t_j \dots t_{i-1} t_i$  is the longest suffix of  $T^i$  that occurs at least twice in  $T^i$ . Therefore,  $s_j t_i$  is the active point of

41

## The Algorithm

create  $\text{STree}(\varepsilon)$

$s \leftarrow \text{root}$

$k \leftarrow 1$

for  $i \leftarrow 1$  to  $n$  do

$(s,k) \leftarrow \text{update}(s,(k,i))$

$(s,k) \leftarrow \text{canonize}(s,(k,i))$

Transforms  $\text{STree}(T^{i-1})$  into  $\text{STree}(T^i)$ .

Input:  $(s,(k,i))$  s.t.  $(s,(k,i-1))$  is the active point of  $\text{STrie}(T^{i-1})$ .

Output:  $(s',k')$  s.t.  $(s',(k',i-1))$  is the endpoint of  $\text{STrie}(T^{i-1})$ .

Input: a reference pair  $(s,(k,p))$  for some state  $r$ .

Output:  $(s',k')$  s.t.  $(s',(k',p))$  is the canonical reference pair for  $r$ .

2

## update( $s,(k,i)$ )

old-r  $\leftarrow$  root

(endpoint,r)  $\leftarrow$  test-and-split( $s,(k,i-1),t_i$ )

while not endpoint do

    create new state  $r'$ ;  $g'(r,(i,\infty)) \leftarrow r'$

    if old-r  $\neq$  root then  $f'(\text{old-r}) \leftarrow r$

    old-r  $\leftarrow$  r

$(s,k) \leftarrow \text{canonize}(f'(s),(k,i-1))$

    (endpoint,r)  $\leftarrow$  test-and-split( $s,(k,i-1),t_i$ )

if old-r  $\neq$  root then  $f'(\text{old-r}) \leftarrow s$

return  $(s,k)$

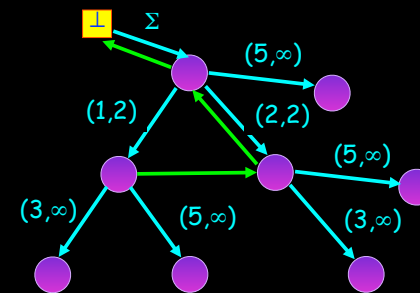
Input: the canonical reference pair for some state  $r$ , and  $t_i$ .

Output: true/false if  $r$  is the endpoint or not, and the explicit state  $r$  (creating it if needed).

43

## update

c a c a o



$s = \perp$

$k = 5$

$i = 5$

44

## test-and-split( $s, (k, p), t$ )

```
if  $k \leq p$  then
  find the  $t_k$ -transition  $g'(s, (k, p')) = s'$  from  $s$ 
  if  $t = t_{k'+p-k+1}$  then return (true,  $s$ )
  else
    create a new state  $r$ 
    replace  $g'(s, (k, p')) = s'$  by  $g'(s, (k, k'+p-k)) = r$ 
      and  $g'(r, (k'+p-k+1, p')) = s'$ 
    return (false,  $r$ )
else if  $\exists t$ -transition from  $s$  then return (false,  $s$ )
```

## canonize( $s, (k, p)$ )

```
if  $p < k$  then return ( $s, k$ )
else
  find the  $t_k$ -transition  $g'(s, (k, p')) = s'$  from  $s$ 
  while  $p' - k' \leq p - k$  do
     $k \leftarrow k + p' - k' + 1$ 
     $s \leftarrow s'$ 
    if  $k \leq p$  then
      find the  $t_k$ -transition  $g'(s, (k, p')) = s'$  from  $s$ 
  return ( $s, k$ )
```

46

## Running Time

**Theorem:** The running time of the algorithm is  $O(n)$ .

**Proof:** We divide the running time into two components:

1. The total time of the procedure canonize.
2. The rest.

47

## update

```
old-r  $\leftarrow$  root
(endpoint,  $r$ )  $\leftarrow$  test-and-split( $s, (k, i-1), t$ )
while not endpoint do
  create new state  $r'$ ;  $g'(r, (i, \infty)) \leftarrow r'$ 
  if old-r  $\neq$  root then  $f'(\text{old-r}) \leftarrow r$ 
  old-r  $\leftarrow r$ 
  ( $s, k$ )  $\leftarrow$  canonize( $f'(s), (k, i-1)$ )
  (endpoint,  $r$ )  $\leftarrow$  test-and-split( $s, (k, i-1), t$ )
if old-r  $\neq$  root then  $f'(\text{old-r}) \leftarrow s$ 
```

Called  $n$   
times

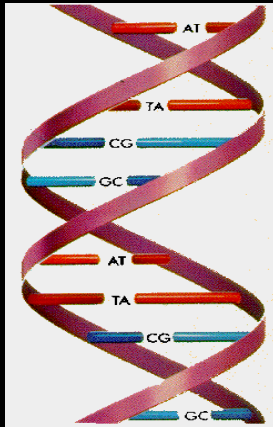
In each  
execution of the  
loop, a new state  
is created.

$O(1)$

48



## Applications in Biology



53

## Finding Repeats in DNA

- The DNA contains many repetitive sequences with different biological functions.
- We want to find all maximal repeats in a DNA sequence.

`ACCAGTTTCGCGCATGAACGTTTCGACCGGTTTCGAT`

54

## Finding Repeats in DNA (cont.)

**Theorem:** All maximal repeats in a sequence  $T$  can be found in  $O(|T|)$  time using suffix trees.

55

## Finding Repeats in DNA (cont.)

**Lemma:** If  $w$  is a maximal repeat in  $T$ , then the state  $w$  in  $S\text{Tree}(T)$  is explicit.

**Proof:** If  $w$  is a maximal repeat then there are at least two occurrences of  $w$  in  $T$  s.t. the character following  $w$  is different. Thus  $w$  is a branching state, and therefore it is explicit.

56

## Finding Repeats in DNA (cont.)

Corollary: There are at most  $O(|T|)$  maximal repeats in  $T$ .

Proof: By the above lemma, each maximal repeat corresponds to an explicit state. Since  $S\text{Tree}(T)$  has  $O(|T|)$  explicit states,  $T$  has  $O(|T|)$  maximal repeats.

57

## Finding Repeats in DNA (cont.)

Definition: The **left character** of a leaf  $t_{i\dots t_n}$  of  $S\text{Tree}(T)$  is  $t_{i-1}$ .

Definition: A node  $w$  of  $S\text{Tree}(T)$  is called **left diverse** if there are at least two leaves in  $w$ 's subtree with different left characters.

Note that, by definition, a left diverse node is not a leaf.

58

## Finding Repeats in DNA (cont.)

Lemma: A substring  $w$  of  $T$  is a maximal repeat iff  $w$  is a left diverse explicit state in  $S\text{Tree}(T)$ .

59

## Finding Repeats in DNA (cont.)

Proof:

1. Suppose  $w$  is a maximal repeat.
  - i. By the previous lemma  $w$  is explicit.
  - ii.  $\exists a \neq b \in \Sigma$  s.t.  $aw$  and  $bw$  are substrings of  $T$ . Let  $awu$  and  $bwv$  be the corresponding suffixes.  
 $\Rightarrow wu$  and  $wv$  are two leaves in the subtree of  $w$  with different left characters.

60

## Finding Repeats in DNA (cont.)

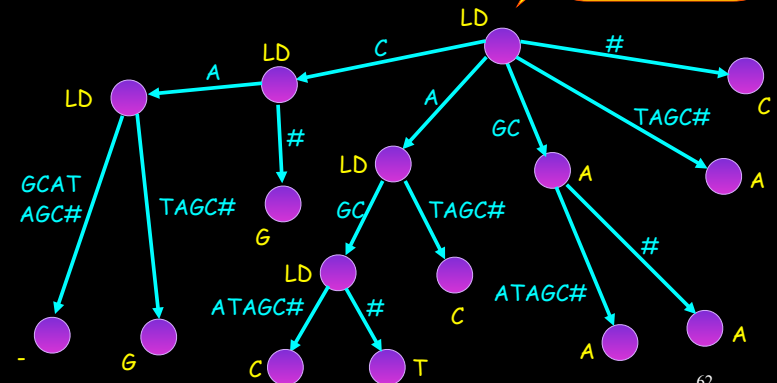
2. Suppose that  $w$  is explicit and left diverse.

	$aw$	$bw$
(i)	$awc$	$bwd$
(ii)	$awc$	$bwc$
		$wd$

61

## Finding Repeats in DNA (cont.)

$CAGCATAGC$



62

## Bibliography

- On-Line Construction of Suffix Trees  
E. Ukkonen
- Algorithms on String, Trees, and Sequences  
Dan Gusfield

63