

---

# Advanced Data Structures

String Pattern Matching/Text Search

# What is Pattern Matching?

---

## Definition:

given a text string  $T$  and a pattern string  $P$ ,  
find the pattern inside the text

$T$ : the rain in spain stays mainly on the plain

$P$ : nth

# Text search

---

Pattern matching directly

- Brute force

- KMP

- BM

Regular expressions (Not in this course)

Indices for pattern matching

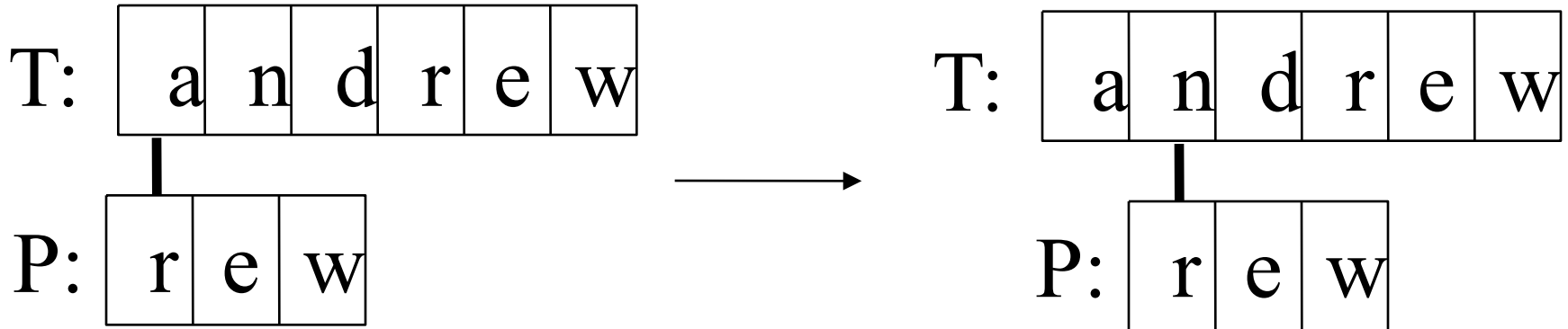
- Inverted files

and

# The Brute Force Algorithm

---

Check each position in the text T to see if the pattern P starts in that position



P moves 1 char at a time through T



....



---

The brute force algorithm is fast when the alphabet of the text is large

e.g. A..Z, a..z, 1..9, etc.

It is slower when the alphabet is small

e.g. 0, 1 (as in binary files, image files, etc.)

*continued*

---

E ample of a worst case:

T: "aaaaaaaaaaaaaaaaaaaaaaaaah"

P: "aaah"

E ample of a more average case:

T: "a string searching e ample is standard"

P: "store"

# The KMP Algorithm

---

The Knuth-Morris-Pratt (KMP) algorithm looks for the pattern in the text in a *left-to-right* order (like the brute force algorithm).

But it shifts the pattern more intelligently than the brute force algorithm.

*continued*

# Summar

---

If a mismatch occurs between the text and pattern  $P$  at  $P[j]$ , what is the *most* we can shift the pattern to avoid wasteful comparisons?

# Summar

---

If a mismatch occurs between the text and pattern  $P$  at  $P[j]$ , what is the *most* we can shift the pattern to avoid wasteful comparisons?

*Answer:* the largest prefix of  $P[0 \dots j-1]$  that is a suffix of  $P[1 \dots j-1]$

[illegible]

# KMP Advantages

---

KMP runs in optimal time:  $O(m+n)$   
very fast

The algorithm never needs to move  
backwards in the input text,  $T$

this makes the algorithm good for processing  
very large files that are read in from external  
devices or through a network stream

# KMP Disadvantages

---

KMP doesn't work so well as the size of the alphabet increases

- more chance of a mismatch (more possible mismatches)

- mismatches tend to occur early in the pattern, but KMP is faster when the mismatches occur later

A fast string searching algorithm. *Communications of the ACM*.  
Vol. 20 p.p. 762-772, 1977.

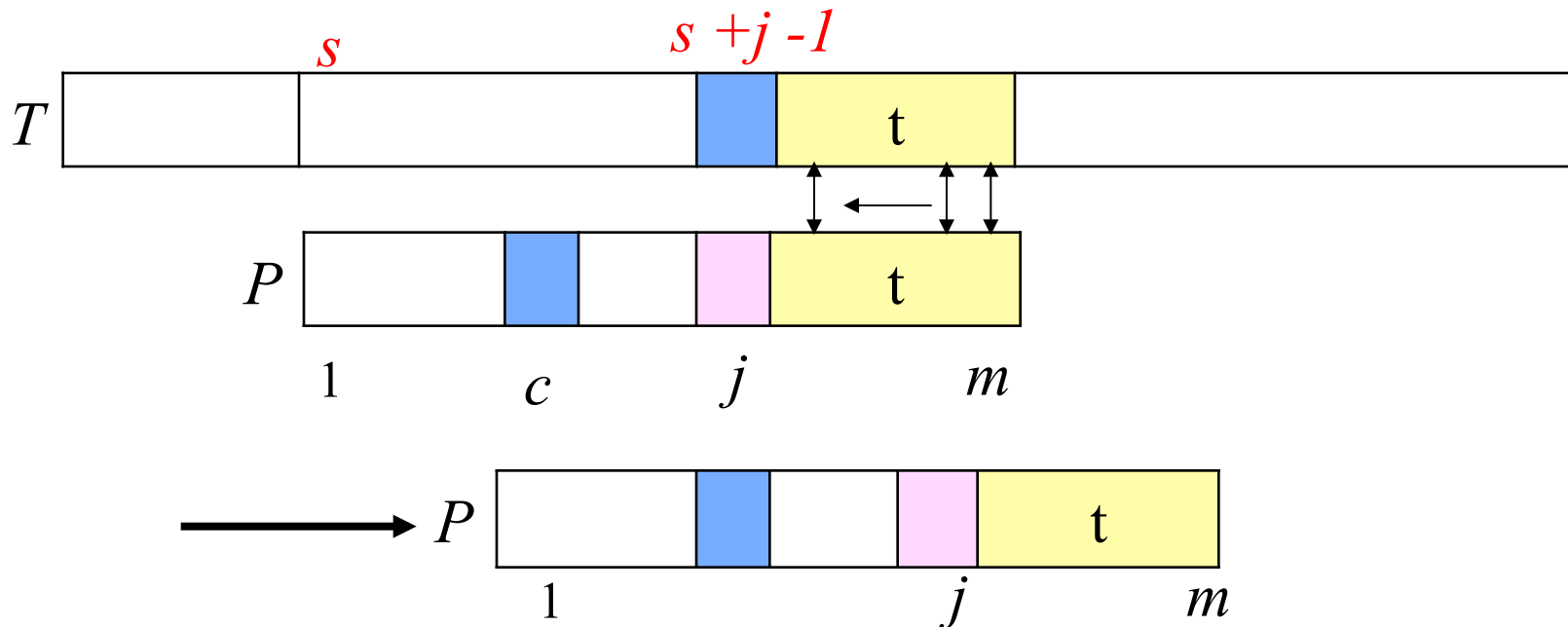
, . . and , . .

The algorithm compares the pattern  $P$  with the substring of sequence  $T$  within a sliding window in the - - .

The and  
are used to determine the movement of sliding window.

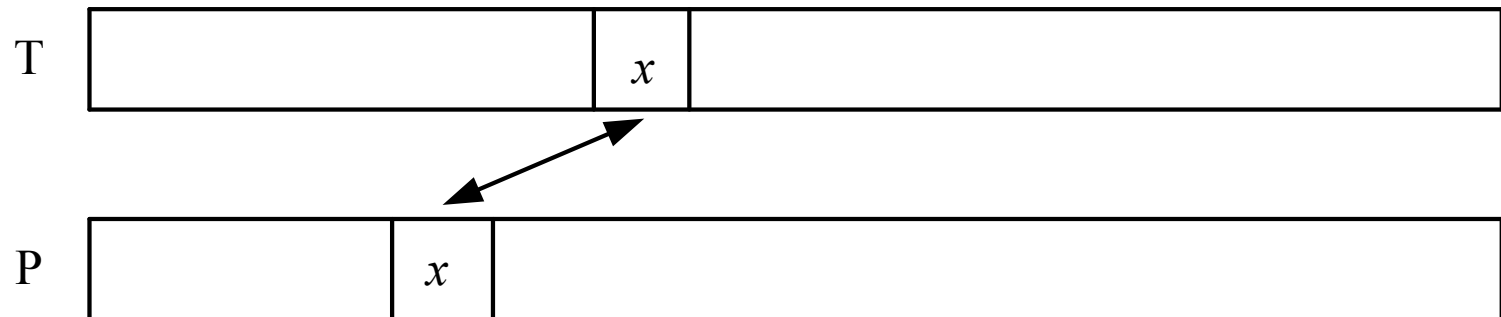
Suppose that  $P_l$  is aligned to  $T_s$  now, and we perform a pairwise comparing between the  $T$  and pattern  $P$  from right to left. Assume that the first mismatch occurs when comparing  $T_{s+j-l}$  with  $P_j$ .

Since  $T_{s+j-l} \neq P_j$ , we move the pattern  $P$  to the right such that the largest position  $c$  in the left of  $P_j$  is equal to  $T_{s+j-l}$ . We can shift the pattern at least  $(j-c)$  positions right.

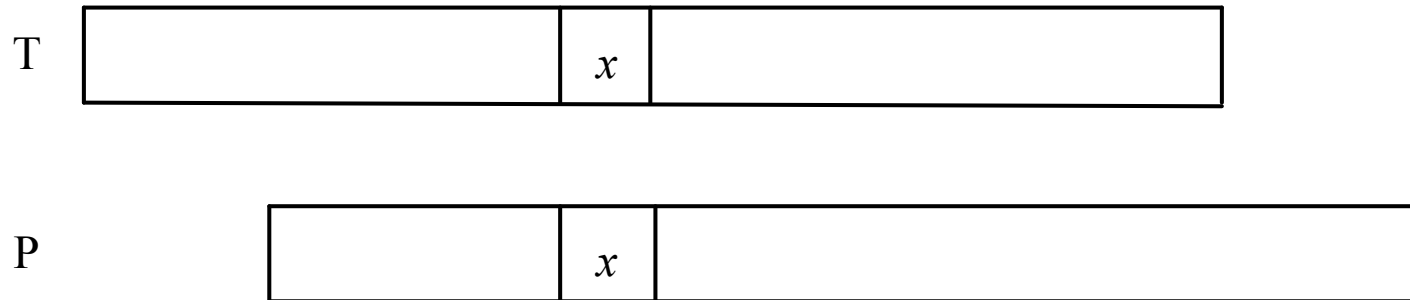


Bad character rule uses Rule 2-1 (Character Matching Rule).

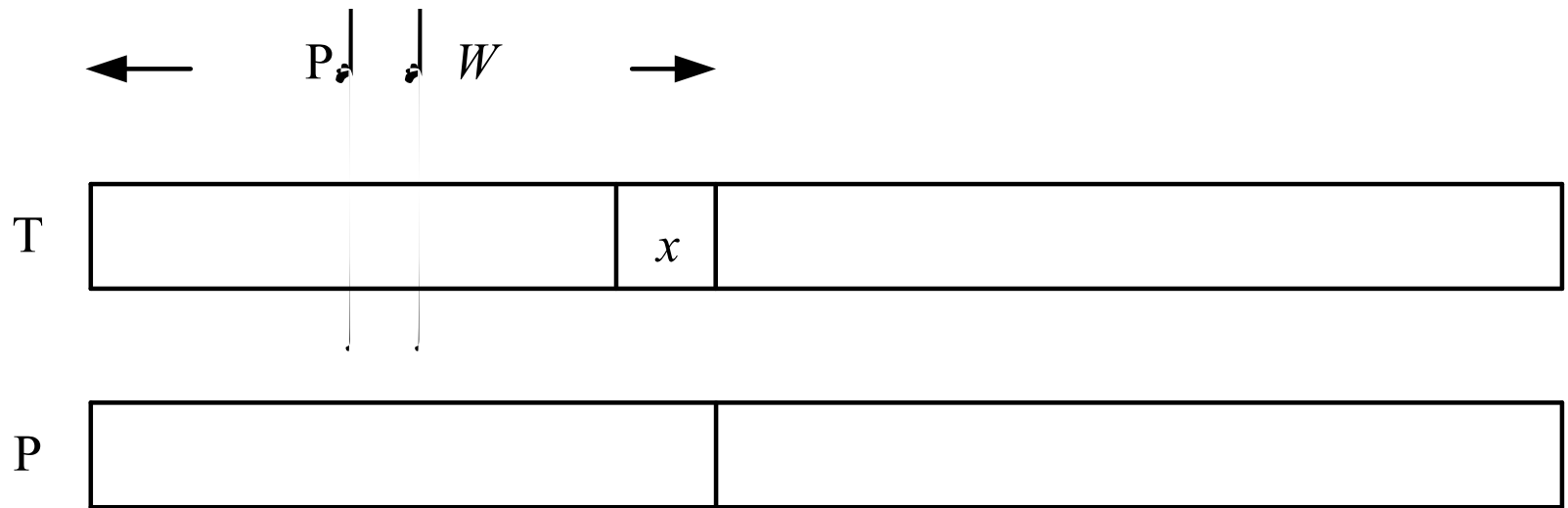
For an character  $x$  in  $T$ , find the nearest  $x$  in  $P$  which is to the left of  $x$  in  $T$ .



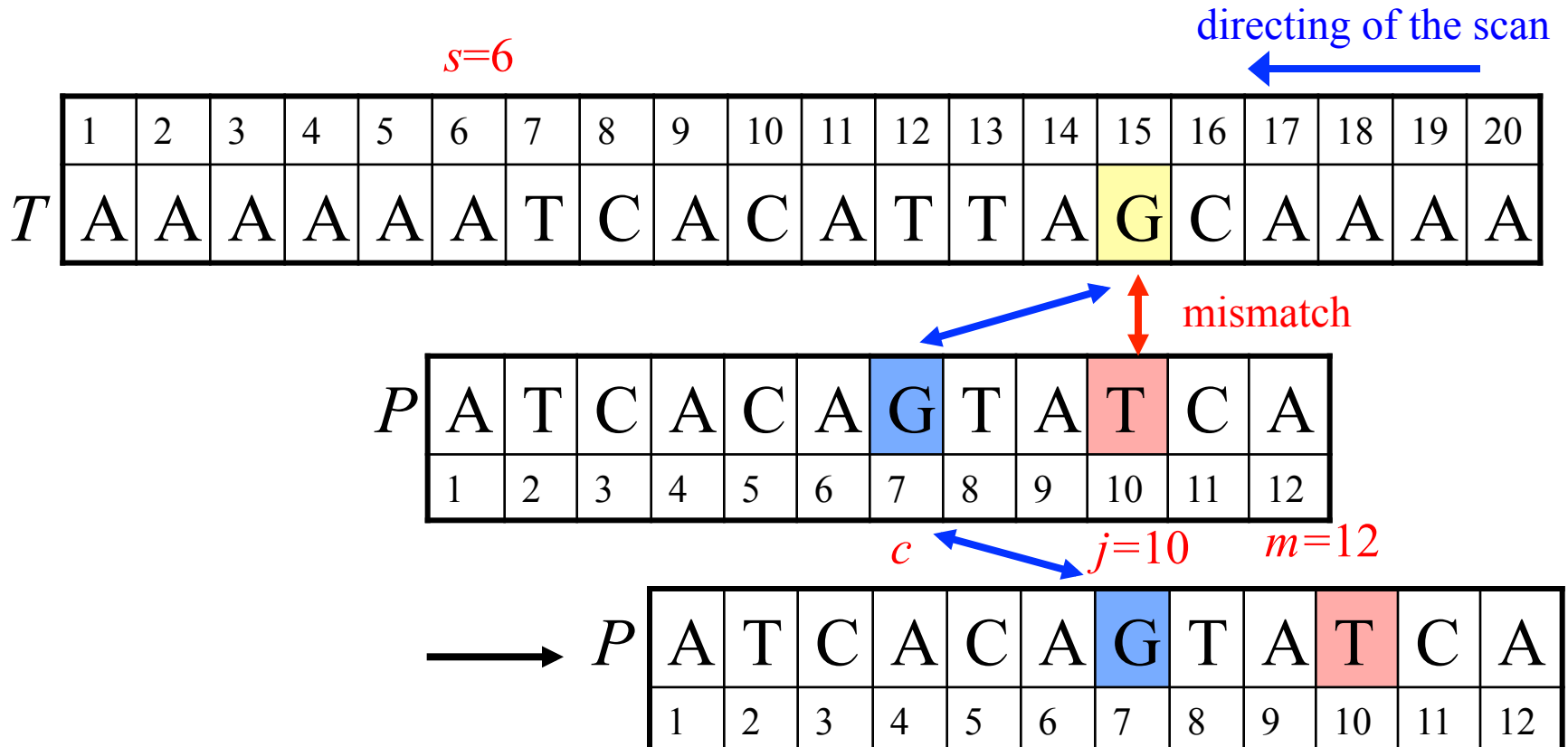
Case 1. If there is a  $x$  in  $P$  to the left of  $T$ , move  $P$  so that the two  $s$  match.



Case 2: If no such a  $x$  exists in  $P$ , consider the partial window defined by  $x$  in  $T$  and the string to the left of it.

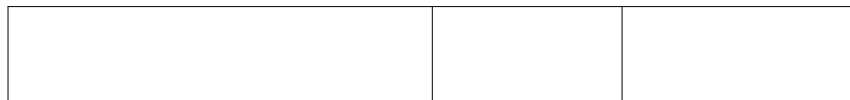
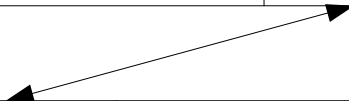
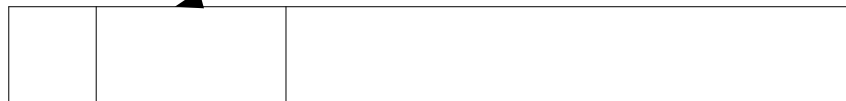


E : Suppose that  $P_l$  is aligned to  $T_6$  now. We compare pairwise between  $T$  and  $P$  from right to left. Since  $T_{16,17} = P_{11,12} =$  “CA” and  $T_{15} =$  “G”  $\neq P_{10} =$  “T”. Therefore, we find the rightmost position  $c=7$  in the left of  $P_{10}$  in  $P$  such that  $P_c$  is equal to “G” and we can move the window at least  $(10-7=3)$  positions.

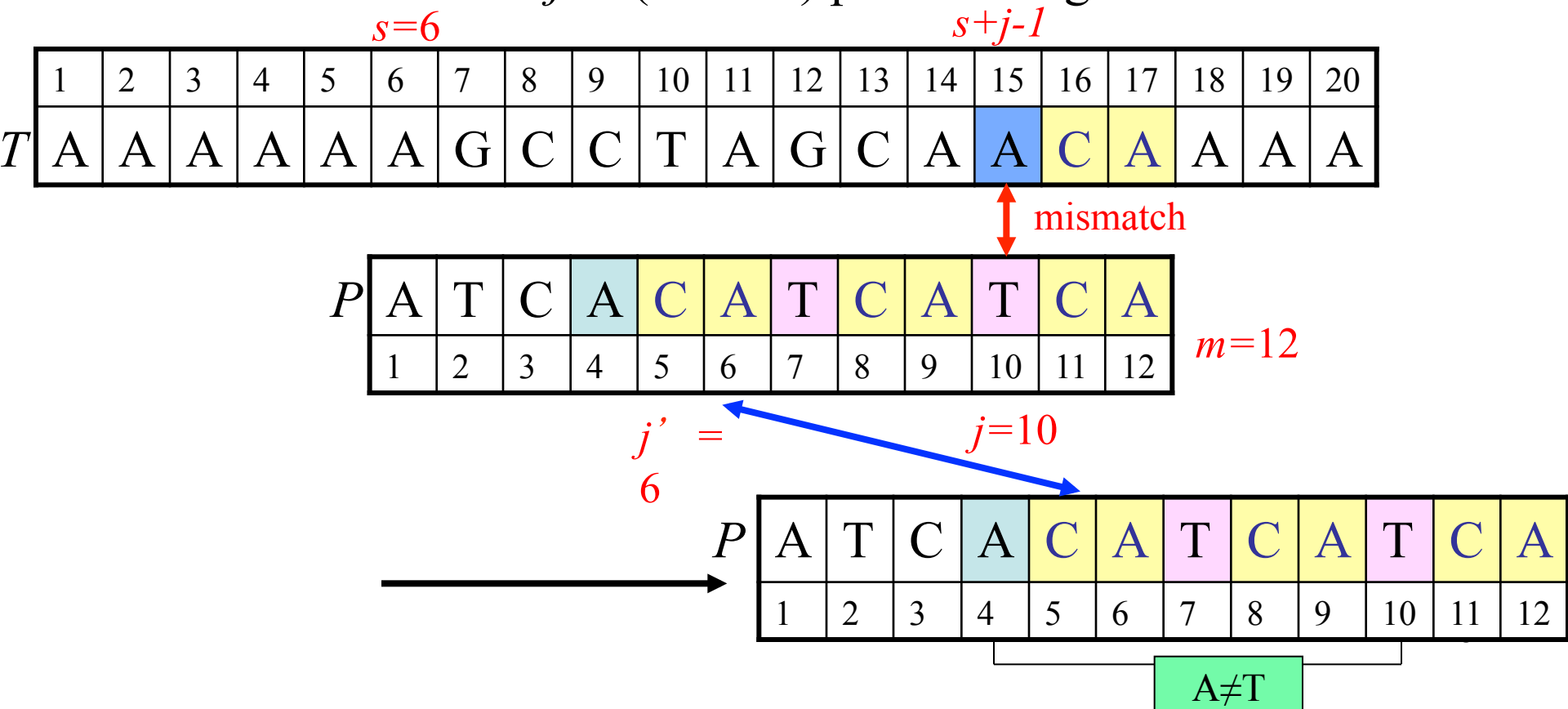


If a mismatch occurs in  $T_{s+j-1}$ , we match  $T_{s+j-1}$  with  $P_{j'-m+j}$ , where  $j'$  ( $m-j+1 \leq j' < m$ ) is the

For an substring  $u$  in  $T$ , find a nearest  $u$  in  $P$  which is to the left of it. If such a  $u$  in  $P$  exists, move  $P$ ; otherwise, we may define a new partial window.

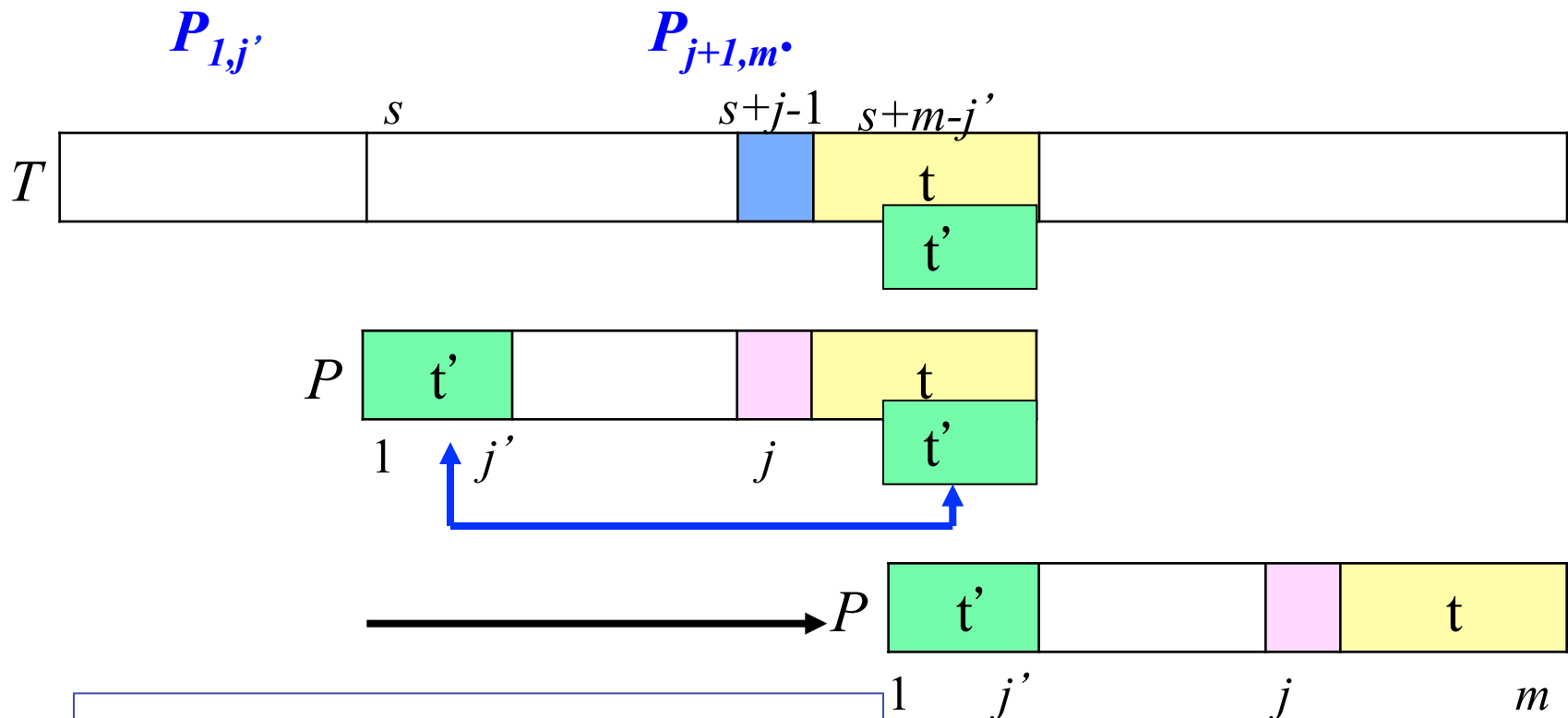


E : Suppose that  $P_l$  is aligned to  $T_6$  now. We compare pairwise between  $P$  and  $T$  from right to left. Since  $T_{16,17} = \text{"CA"} = P_{11,12}$  and  $T_{15} = \text{"A"} \neq P_{10} = \text{"T"} . We find the substring \text{"CA"} in the left of  $P_{10}$  in  $P$  such that \text{"CA"} is the suffix of  $P_{1,6}$  and the left character to this substring \text{"CA"} in  $P$  is not equal to  $P_{10} = \text{"T"} . Therefore, we can move the window at least  $m-j'$  (12-6=6) positions right.$$



Good Suffi Rule 2 is used only when Good Suffi Rule 1 can not be used. That is,  $t$  does not appear in  $P(1, j)$ . Thus,  $t$  is in  $P$ .

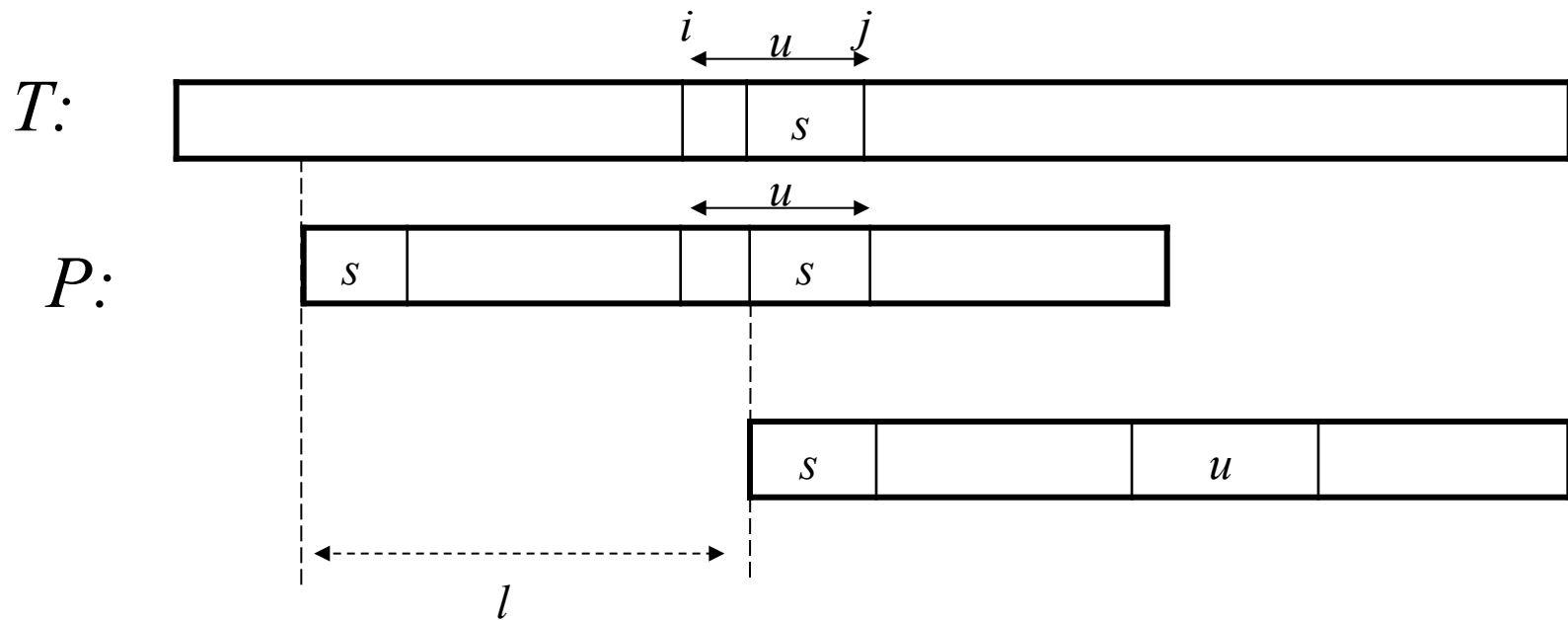
If a mismatch occurs in  $T_{s+j-1}$ , we match  $T_{s+m-j'}$  with  $P_1$ , where  $j'$  ( $1 \leq j' \leq m-j$ ) is such that



P.S. :  $t'$  is suffi of substring  $t$ .

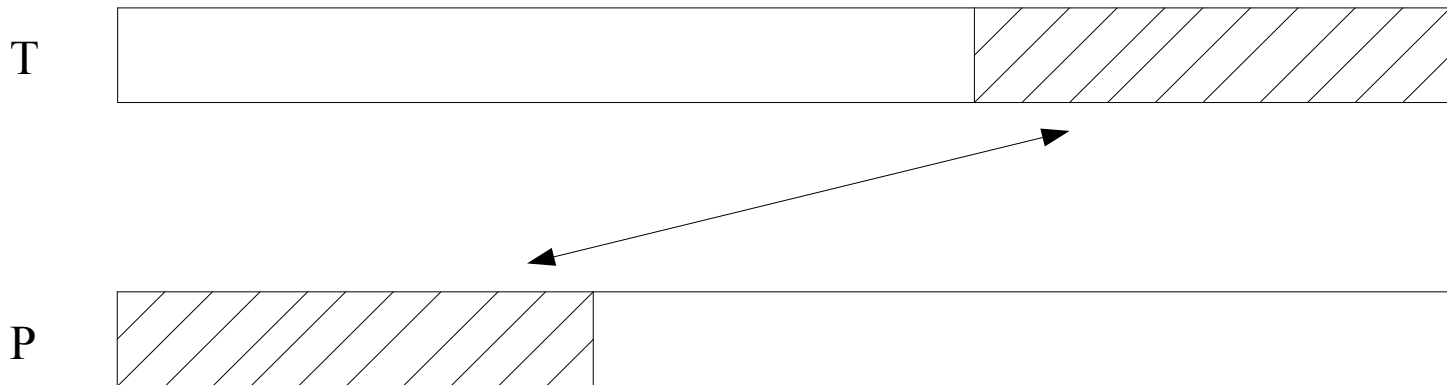
The substring  $u$  appears in  $P$  exactly once.

If the substring  $u$  matches with  $T_{i,j}$ , no matter whether a mismatch occurs in some position of  $P$  or not, we can slide the window by  $l$ .



The string  $s$  is the longest prefix of  $P$  which equals to a suffix of  $u$ .

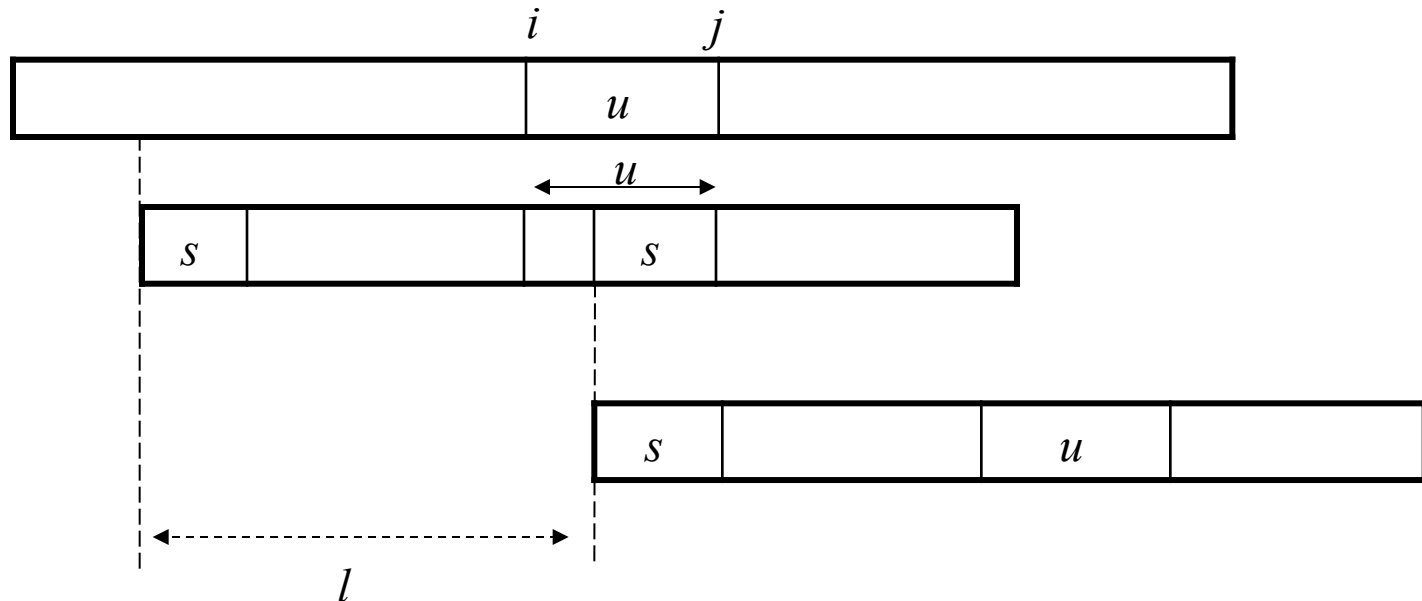
For a window to have an chance to match a pattern, in some wa , there must be a suffi of the window which is equal to a pref of the pattern.



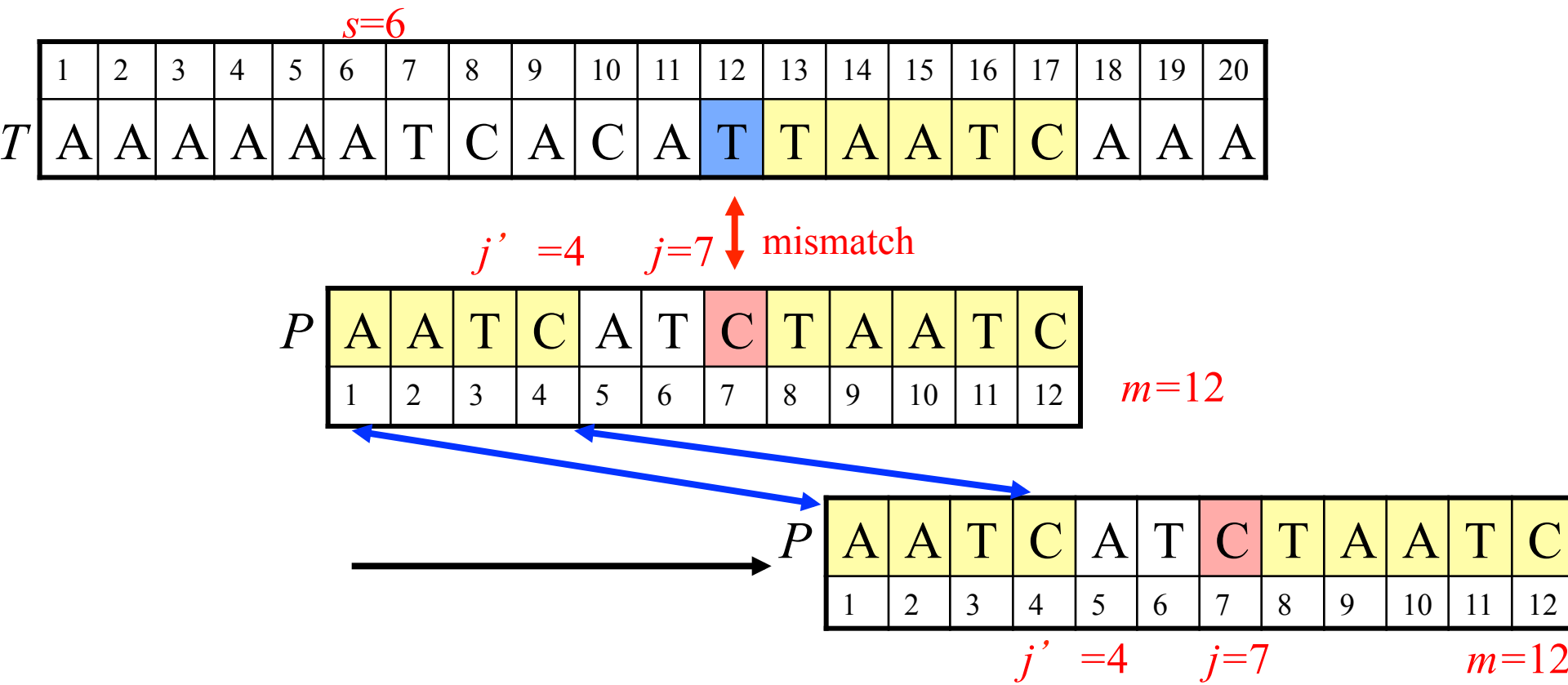
Note that the above rule also uses Rule 1.

It should also be noted that the unique substring is the shorter and the more right-sided the better.

A short  $u$  guarantees a short (or even empty)  $s$  which is desirable.



E : Suppose that  $P_1$  is aligned to  $T_6$  now. We compare pair-wise between  $P$  and  $T$  from right to left. Since  $T_{12} \neq P_7$  and there is no substring  $P_{8,12}$  in left of  $P_8$  to exactly match  $T_{13,17}$ . We find a longest suffix “AATC” of substring  $T_{13,17}$ , the longest suffix is also prefix of  $P$ . We shift the window such that the last character of prefix substring to match the last character of the suffix substring. Therefore, we can shift at least  $12-4=8$  positions.



Let  $Bc(a)$  be the rightmost position of  $a$  in  $P$ . The function will be used for applying *bad character rule*.

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A

$\Sigma$	A	C	G	T
$B$	12	11	0	10

We can move our pattern right  $j - B(T_{s+j-})$  position based on  $Bc$  function.

$j$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$T$	A	G	C	T	A	G	C	C	T	G	C	A	C	G	T	A	C	A

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A

Move  
 $10 - B(G) = 10$  positions

Let  $Gs(j)$  be b *good*  
*suffix rule* when a mismatch occurs for comparing  $P_j$   
with some character in  $T$ .

$gs_1(j)$  be the largest  $k$  such that  $P_{j+1,m} \neq P_{1,k}$ , where  $m-j+1 \leq k < m$ ; 0 if there is no such  $k$ .  
 ( $gs_1$  is for Good Suffi Rule 1)

$gs_2(j)$  be the largest  $k$  such that  $P_{1,k} \neq P_{j+1,m}$ , where  $1 \leq k \leq m-j$ ; 0 if there is no such  $k$ .  
 ( $gs_2$  is for Good Suffi Rule 2.)

$Gs(j)$   $m$   $gs_1, gs_2$ , if  $j = m$ ,  $Gs(j)=1$ .

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$												
$gs_1$	0	0	0	0	0	0	9	0	0	6	1	0
$gs_2$	4	4	4	4	4	4	4	4	1	1	1	0
$Gs$	8	8	8	8	8	8	3	8	11	6	11	1

$$gs_1(7)=9$$

$\therefore P_{8,12}$  is a suffi of  $P_{1,9}$   
 and  $P_4 \neq P_7$

$$gs_2(7)=4$$

$\therefore P_{1,4}$  is a suffi of  $P_{8,12}$

How do we obtain  $gs_1$  and  $gs_2$ ?

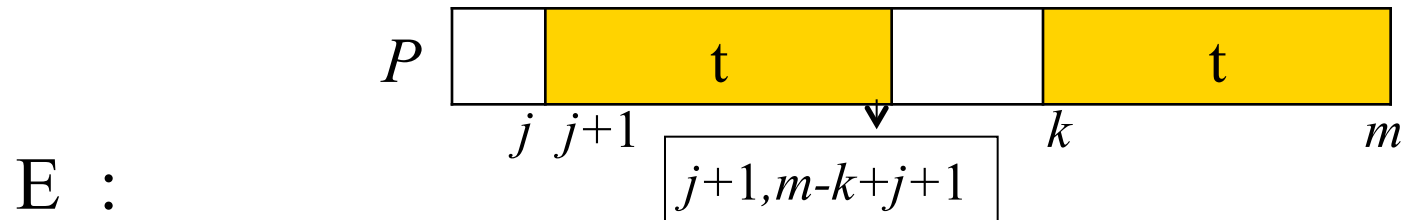
In the following, we shall show that by  
constructing the  $gs_1$  and  $gs_2$ , we can kill  
two birds with one arrow.

$f'$

For  $1 \leq j \leq m-1$ , let the suffix function  $f'(j)$  for  $P_j$  be the  $k$  such that  $P_{k,m} = P_{j+1,m-k+j+1}$ ; ( $j+1 \leq k \leq m$ )

If there is no such  $k$ , we set  $f' = m+1$ .

If  $j=m$ , we set  $f'(m)=m+2$ .



$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14

- (4)=8, it means that  $P_{f'(4),m} = P_{8,12} = P_{5,9} = P_{4+1,4+1+m-f'(4)}$   
 Since there is no  $k$  for  $13=j+2 \leq k \leq 12$ , we set  $f'(11)=13$ .

Suppose that the Suffi is obtained. How can we use it to obtain  $gs_1$  and  $gs_2$ ?

$gs_1$  can be obtained by scanning the Suffi function from right to left.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
$T$	G	A	T	C	G	A	T	C	A	A	T	C	A	T	C	A	C	A	T	G	A	T	C	A

$P$	A	T	C	A	C	A	T	C	A	T	C	A
	1	2	3	4	5	6	7	8	9	10	11	12

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14

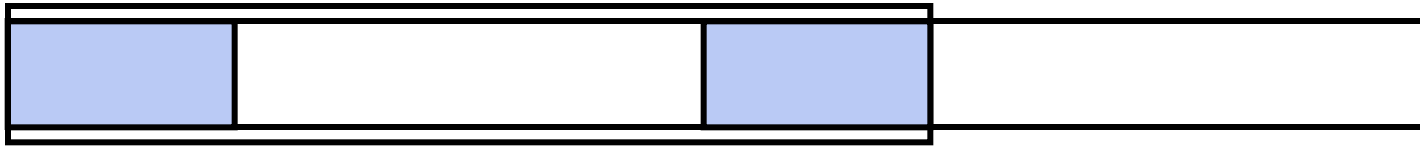
As for Good Suffi Rule 2, it is relative easier.

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14

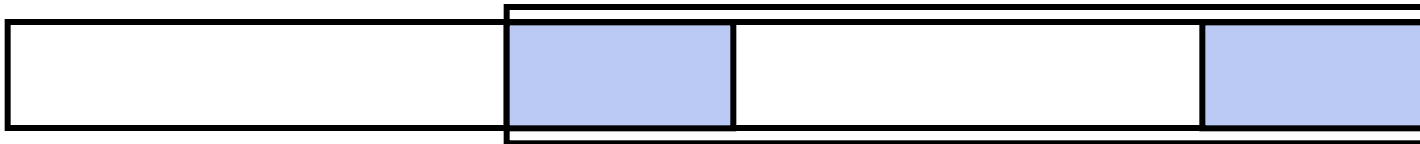
Question: How can we construct the Suffix function?

To explain this, let us go back to the prefix function used in the KMP Algorithm.

The following figure illustrates the prefix function in the KMP Algorithm.



The following figure illustrates the suffix function of the BM Algorithm.



We now can see that actually the suffix function is the same as the prefix . The only difference is now we consider a suffix . Thus, the recursive formula for the prefix function in KMP Algorithm can be slightly modified for the suffix function in BM Algorithm.

The formula of suffi function  $f'$  as follows :

$$L \quad f'^x(y) = f'(f'^{x-1}(y)) \quad x > 1 \quad f'^1(y) = f'(y)$$

$$f'(j) = \begin{cases} m+2, & j = m \\ f'^k(j+1)-1, & 1 \leq j \leq m-1 \\ m+1, & \end{cases} \quad k \geq 1 \quad P_{j+1} = P_{f'^k(j+1)-1};$$

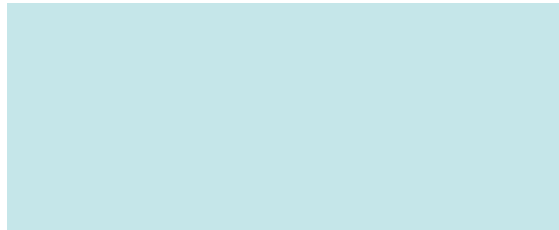
$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f$												14

$$f^j + ,$$

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f$											13	14

$$f^k P_{m+}^+ P_f^{k(j+)-} ,$$

$$k=1 \rightarrow P_{12} \neq P_{13}$$



$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$								12	13	13	13	14

$$\because P_{j+1} = P_{f' (j+1)-1} \Rightarrow P_9 = P_{12},$$

$$f' \quad f' \quad (j+1) - \quad -$$

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$							11	12	13	13	13	14

$$\because P_{j+1} = P_{f' (j+1)-1} \Rightarrow P_8 = P_{11},$$

$$f' \quad f' \quad (j+1) - \quad -$$

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$				8	9	10	11	12	13	13	13	14

$$\because P_{j+1} = P_{f' \quad (j+1)-1} \Rightarrow P_5 = P_8,$$

$$f' \quad (j+1) - \quad -$$

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$			12	8	9	10	11	12	13	13	13	14

$$\because P_{j+1} = P_{f' \quad (j+1)-1} \Rightarrow P_4 = P_{f' \quad (4)-1} = P_{12},$$

$$f' \quad (j+1) - \quad -$$

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$		11	12	8	9	10	11	12	13	13	13	14

$$\because P_{j+1} = P_{f' \quad (j+1)-1} \Rightarrow P_3 = P_{f' \quad (3)-1} = P_{11},$$

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14

$$\because P_{j+1} = P_{f' \quad (j+1)-1} \Rightarrow P_2 = P_{f' \quad (2)-1} = P_{10},$$

Let  $G'(j)$ ,  $1 \leq j \leq m$ , to be the largest number of shifts by good suffix rules.

First, we set  $G'(j)$  to zeros as their initializations.

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14
$G'$	0	0	0	0	0	0	0	0	0	0	0	0

We scan from right to left and  $gs_1(j)$  is determined during the scanning, then  $gs_1(j) \geq gs_2(j)$

If  $P_j = P_4 \neq P_7 = P_{f'(j)-1}$ , we know  $gs(f'(j)-1) = m + j - f'(j) + 1 = 9$ .

- When  $j=12$ ,  $t=13$ .  $t > m$ .
- When  $j=11$ ,  $t=12$ . Since  $P_{11} = 'C' \neq 'A' = P_{12}$ ,  
 $G'(t) = m - gs_1(t), gs_2(t) = m - \underline{gs_1(t)}$   
 $= f'(j) - 1 - j$   
 $\Rightarrow G'(12) = 13 - 1 - 11 = 1$ .

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14
$G'$	0	0	0	0	0	0	0	0	0	0	0	1

$$t = f'(j) - 1 \leq m \quad P_j \neq P_t, G'(t) = f'(j) - 1 - j.$$

$$f^{(k)}(x) = f^{(k-1)}(f'(x) - 1), k \geq 1$$

- When  $j=10, t=12$ . Since  $P_{10} = 'T' \neq 'A' = P_{12}, G'(12) \neq 0$ .
- When  $j=9, t=12$ .  $P_9 = 'A' = P_{12}$ .
- When  $j=8, t=11$ .  $P_8 = 'C' = P_{11}$ .
- When  $j=7, t=10$ .  $P_7 = 'T' = P_{10}$ .
- When  $j=6, t=9$ .  $P_6 = 'A' = P_9$ .
- When  $j=5, t=8$ .  $P_5 = 'C' = P_8$ .
- When  $j=4, t=7$ . Since  $P_4 = 'A' \neq P_7 = 'T', G'(7) = 8 - 1 - 4 = 3$

Besides,  $t = f'^{(2)}(4) - 1 = f(f'(4) - 1) - 1 = 10$ . Since  $P_4 = 'A' \neq P_{10} = 'T', G'(10) = f'(7) - 1 - j = 11 - 1 - 4 = 6$ .

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14
$G'$	0	0	0	0	0	0	3	0	0	6	0	1

$$t \leq m \quad P_j \neq P_{j-1}, G'(t) = f'(j) - 1 - j.$$

$$f^{(k)}(x) = f^{(k-1)}(f'(x)), k \geq 1$$

- When  $j=3$ ,  $t=11$ .  $P_3 = \text{'C'} = P_{11}$ .
- When  $j=2$ ,  $t=10$ .  $P_2 = \text{'T'} = P_{10}$ .
- When  $j=1$ ,  $t=9$ .  $P_1 = \text{'A'} = P_9$ .

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14
$G'$	0	0	0	0	0	0	3	0	6	0	0	1

Based on the above discussion, we can obtain the values using the Good Suffix Rule 1 by scanning the pattern from right to left.

Continuously, we will try to obtain the values using **Good Suffix Rule 2** and those values are still zeros now and scan from left to right.

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14
$G'$	0	0	0	0	0	0	3	0	0	6	0	1

Let  $k'$  be the  $k$  in  $1, \dots, m$  such that  $P_{f'(k)}(1) = P_{f'(k')}(1)$  and  $f'(k)(1) - 1 \leq m$ .

If  $G'(j)$  is not determined in the first scan and  $1 \leq j \leq f'(k')$ , thus, in the second scan, we set  $G'(j) = m - \max\{gs_1(j), gs_2(j)\} = m - gs_2(j) = f'(k')(1) - 2$ . If no such  $k$  exists, set each undetermined value of  $G$  to  $m$  in the second scan.

$k$   $k'$

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14
$G'$	8	8	8	8	8	8	3	8	0	6	0	1

Let  $z$  be  $f^{(k')}(1)-2$ . Let  $k''$  be the  $k$  such that  $f^{(k)}(z)-1 \leq m$ .

Then we set  $G'(j) = m - gs_2(j) = m - (m - f^{(i)}(z) - 1) = f^{(i)}(z) - 1$ , where  $1 \leq i \leq k''$  and  $f^{(i-1)}(z) < j \leq f^{(i)}(z)-1$  and  $f^{(0)}(z) = z$ .

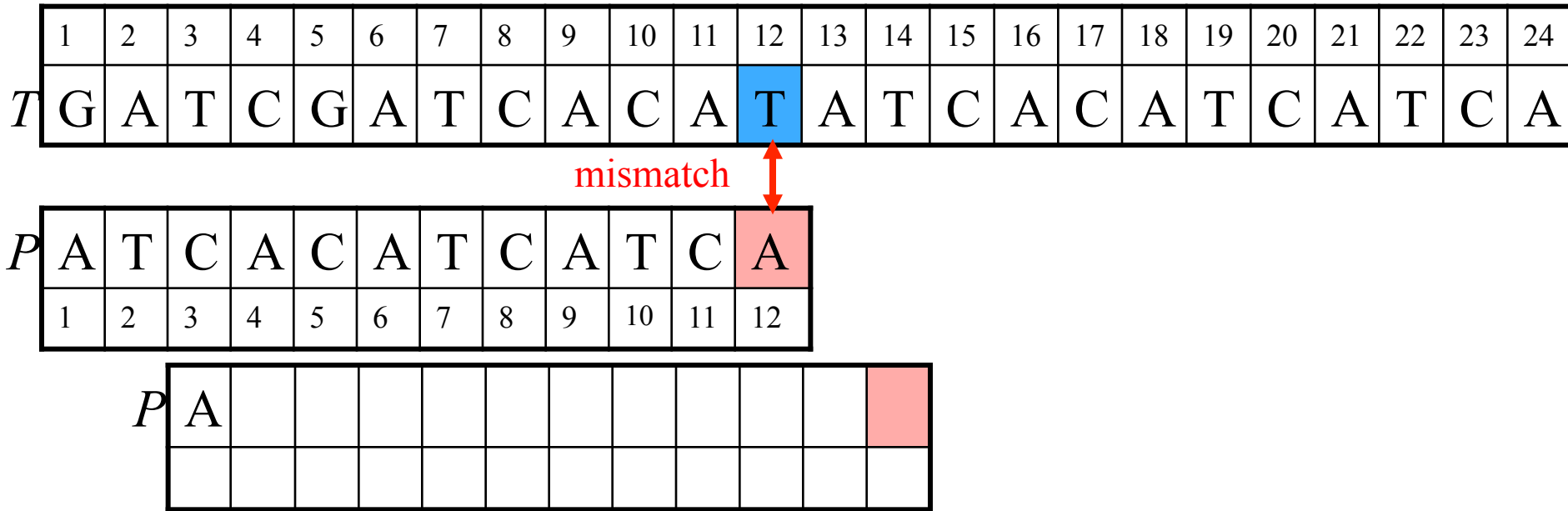
For example,  $z=8$  :

- $k=1, f^{(1)}(8)-1=11 \leq m=12$
- $k=2, f^{(2)}(8)-1=12 \leq m=12 \Rightarrow k''=2$
- $i=1, f^{(0)}(8)-1 = 7 < j \leq f^{(1)}(8)-1=11$ .
- $i=2, f^{(1)}(8)-1 = 11 < j \leq f^{(2)}(8)-1=12$ .
- We set  $G(9)$  and  $G(11)=f^{(1)}(8) - 1 = 12-1 = 11$ .

$j$	1	2	3	4	5	6	7	8	9	10	11	12
$P$	A	T	C	A	C	A	T	C	A	T	C	A
$f'$	10	11	12	8	9	10	11	12	13	13	13	14
$G'$	8	8	8	8	8	8	3	8	11	6	11	1

We essentially have to decide the maximum number of steps.  
 We can move the window right when a mismatch occurs. This is decided by the following function:

$$\text{max } G'(j), j-B(T_{s+j-1})$$



We compare  $T$  and  $P$  from right to left. Since  $T_{12} = \text{"T"} \neq P_{12} = \text{"A"}$ , the largest movement =  $\max_{G'}(j, j - B(T_{s+j-1})) = \max_{G'}(12, 12 - B(T_{12})) = \max_{G'}(1, 12 - 10) = 2$ .

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
<i>T</i>	G	A	T	C	G	A	T	C	A	C	A	T	A	T	C	A	C	A	T	C	A	T	C	A

mismatch



<i>P</i>	A	T	C	A	C	A	T	C	A	T	C	A
	1	2	3	4	5	6	7	8	9	10	11	12

→

<i>P</i>	A	T	C	A	C	A	T	C	A	T	C	A
	1	2	3	4	5	6	7	8	9	10	11	12

<i>j</i>	1	2	3	4	5	6	7	8	9	10	11	12
<i>P</i>	A	T	C	A	C	A	T	C	A	T	C	A
<i>f'</i>	10	11	12	8	9	10	11	12	13	13	13	14
<i>G'</i>	8	8	8	8	8	8	3	8	11	6	11	1

$\Sigma$	A	C	G	T
<i>B</i>	12	11	0	10

After moving, we compare *T* and *P* from right to left. Since  $T_{14} = \text{"T"} \neq P_{12} = \text{"A"}$ ,  
the largest movement =  $\text{ma}_{G'}(j), j-B(T_{s+j-1}) = \text{ma}_{G'}(12), 12-B(T_{14})$   
=  $\text{ma}_{1, 12-10} = 2$ .

The preprocessing phase in  $O(m+\Sigma)$  time and space comple it and searching phase in  $O(mn)$  time comple it .

The worst case time comple it for the ***Boyer-Moore*** method would be  $O((n-m+1)m)$ .

It was proved that this algorithm has  $O(m)$  comparisons when  $P$  is not in  $T$ . However, this algorithm has  $O(mn)$  comparisons when  $P$  is in  $T$ .

, AHO, A.V.,  
 , Volume A , Chapter 5 Elsevier , Amsterdam , 1990, pp. 255-300.  
 , Jun-ichi, A. , IEEE  
 Computer Society Press , 1994.

, BAASE, S. and  
 VAN GELDER, A. , - , Chapter 11 , 1999.

, BAEZA-YATES, R. , NAVARRO, G. and RIBEIRO-  
 NETO, B. , , Chapter 8 , 1999 , pp. 191-228.

,  
 P., Masson Paris , Chapter 10 , 1992 , pp. 337-377.

, BEAUQUIER, D., BERSTEL, J. and CHRISTIENNE,  
 , BOYER R.S. and MOORE J.S. ,  
 , Vol 20 , 1977 , pp. 762-772 .

-  
 , COLE, R., , Vol 23 , 1994 , pp.  
 1075-1091.

, CORMEN, T.H. , LEISERSON, C.E. and RIVEST,  
 R.L. , , Chapter 34 , 1990 , pp. 853-885.

- , CROCHEMORE, M. ,  
 , Chapter 1 , 1997 , pp 1-53

, CROCHEMORE, M. and HANCART, C. ,  
 , Chapter 11, 1999 , pp.  
 11-1-11-28.

LECROQ, T. , , CROCHEMORE, M. and , Chapter 8 ,  
1996 , pp. 162-202.

, CROCHEMORE, M. and RYTTER, W. ,  
1994.

BAEZA-YATES, R.A. , - , GONNET, G.H. and  
251-288,. , Chapter 7 , 1991 , pp.

, GOODRICH, M.T. and TAMASSIA, R. ,  
& , Chapter 11 , 1998 , pp. 441-467.

, GUSFIELD, D. , , 1997.

HANCART, C., Universit Paris 7, France , 1993.

, KNUTH, D.E. , MORRIS, J.H. and PRATT, V.R. ,  
, 1977 , pp.323-350.

LECROQ, T., 1992, *Recherches de mot*., Thesis, Universit of Orl ans, France.

, LECROQ, T. , -  
& , Vol 25 , 1995 , pp. 727-765.

, SEDGEWICK, R. , - , Chapter 19 ,  
1988 , pp. 277-292.

, SEDGEWICK, R. , - ,  
Chapter 19 , 1988.

, STEPHEN, G.A. , World Scientific , 1994.

, WATSON, B.W. ,  
, 1995.

& , WIRTH, N. , Prentice-Hall , Chapter 1 , 1986 , pp.  
17-72.

---

# Suffix trees and suffix arrays

# String/Pattern Matching

- You are given a source string  $S$ .
- Answer queries of the form: is the string  $p_i$  a substring of  $S$ ?
- Knuth-Morris-Pratt (KMP) string matching.
  - $O(|S| + |p_i|)$  time per query.
  - $O(n|S| + \sum |p_i|)$  time for  $n$  queries.
- Suffix tree solution.
  - $O(|S| + \sum |p_i|)$  time for  $n$  queries.

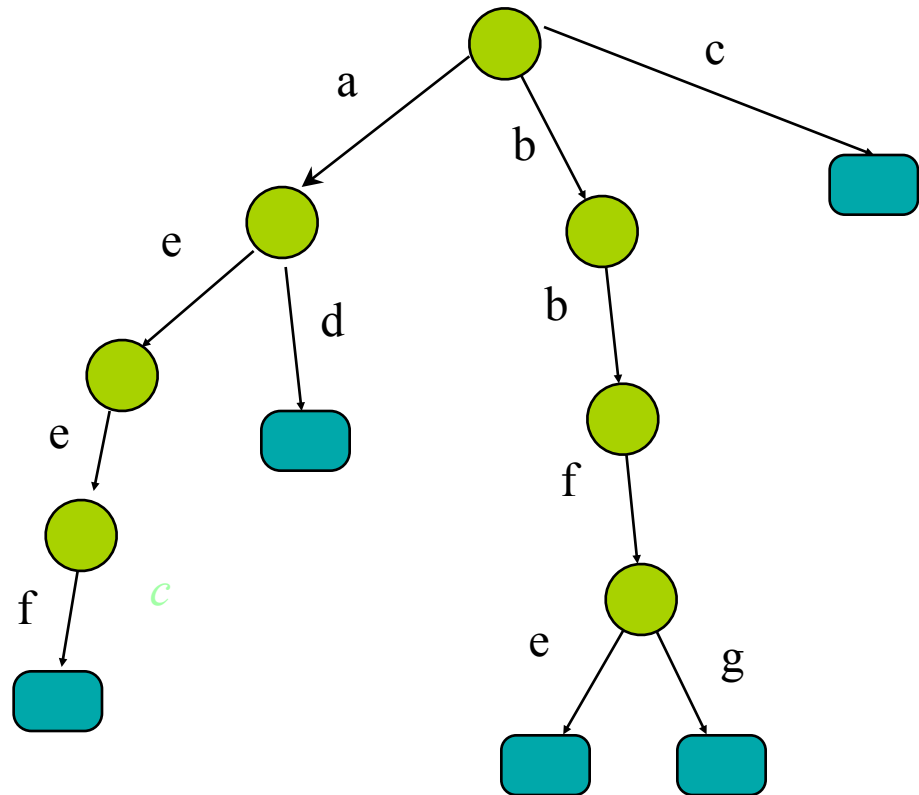
# String/Pattern Matching

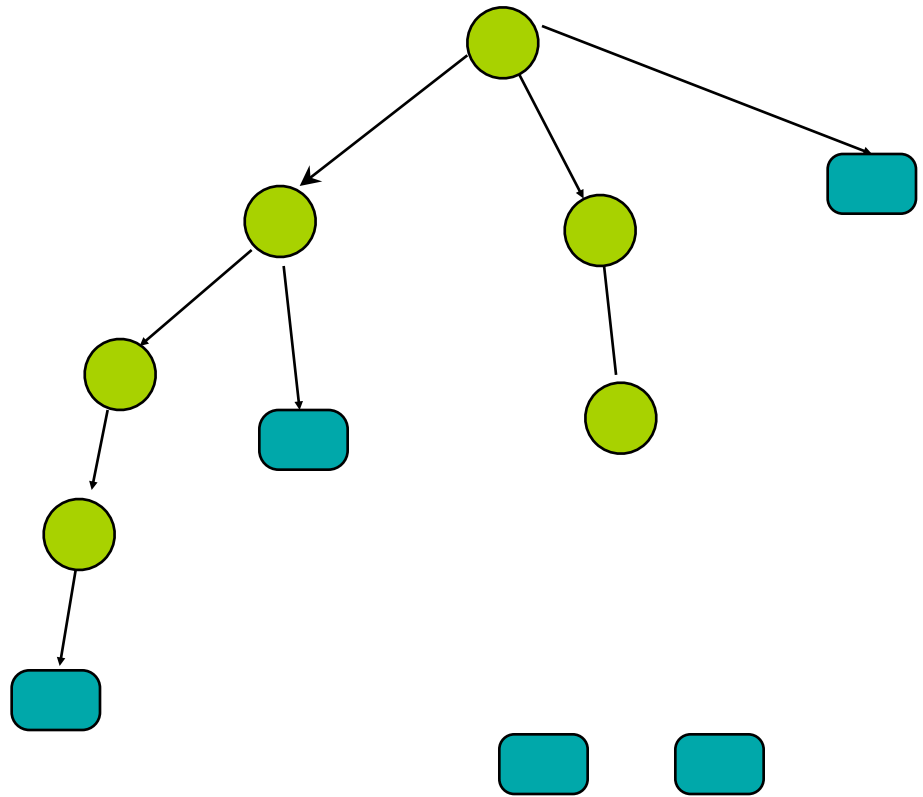
- KMP/BM preprocesses the query string  $p_i$ , whereas the suffix tree method preprocesses the source string  $S$ .

# Trie

A tree representing a set of strings.

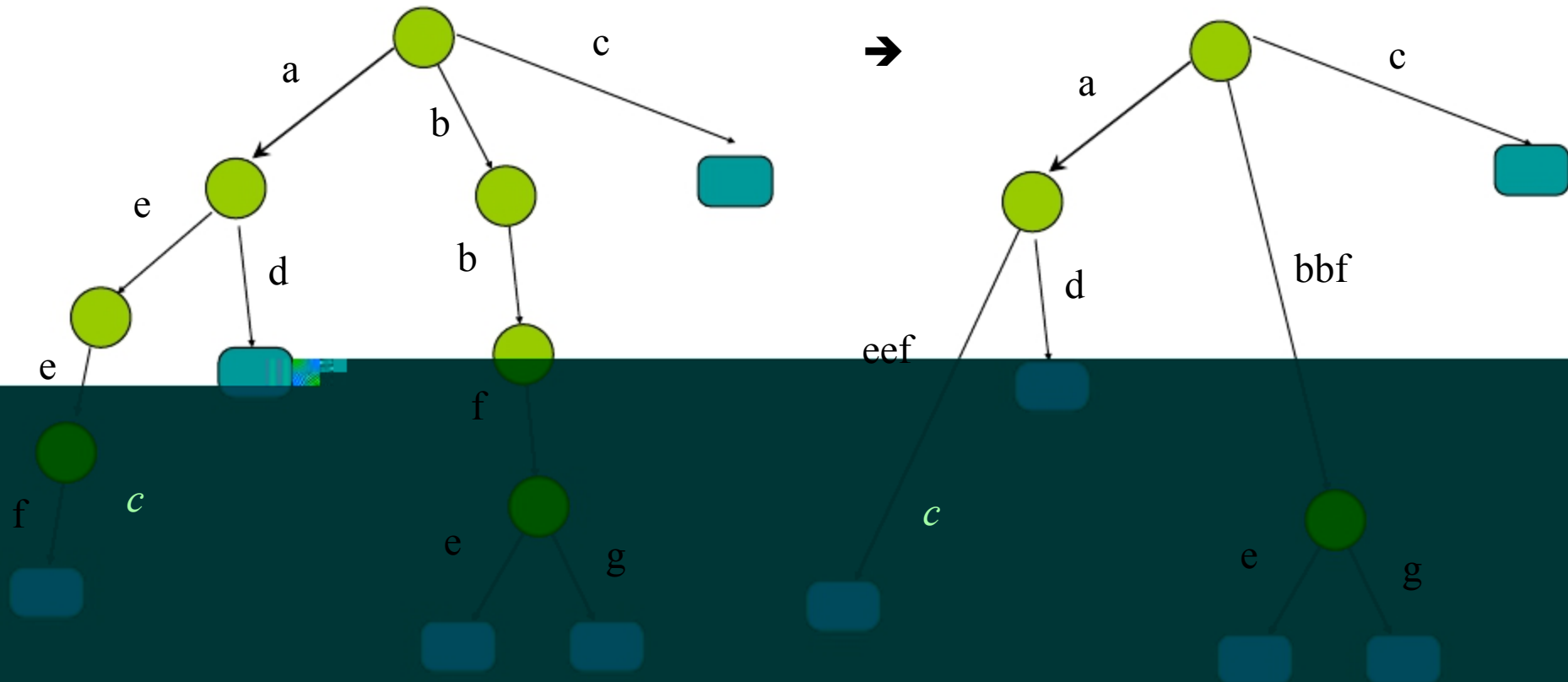
aeef  
ad  
bbfe  
bbfg  
c





# Compressed Trie

Compress unar nodes, label edges b strings



# Suffix tree

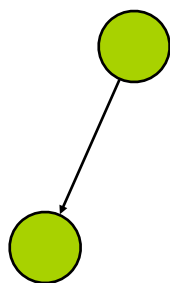
---

Given a string **s** a suffix tree of **s** is a compressed trie of all suffixes of s

To make these suffixes prefix-free we

# The suffix tree $\text{Tree}(T)$ of $T$

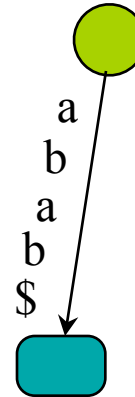
- data structure **suffix tree**,  $\text{Tree}(T)$ , is compacted trie that represents all the suffixes of string  $T$
- linear size:  $\text{Tree}(T) = O(|T|)$
- can be constructed in linear time  $O(|T|)$
- has *myriad virtues* (A. Apostolico)
- is well-known: Google hits



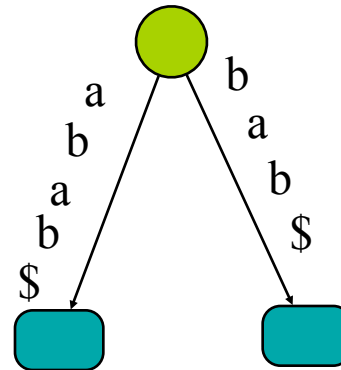
# Trivial algorithm to build a Suffix tree

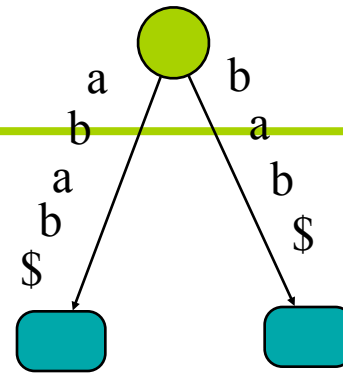
---

Put the largest suffix in

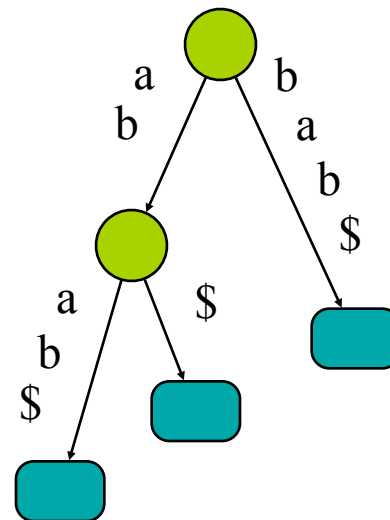


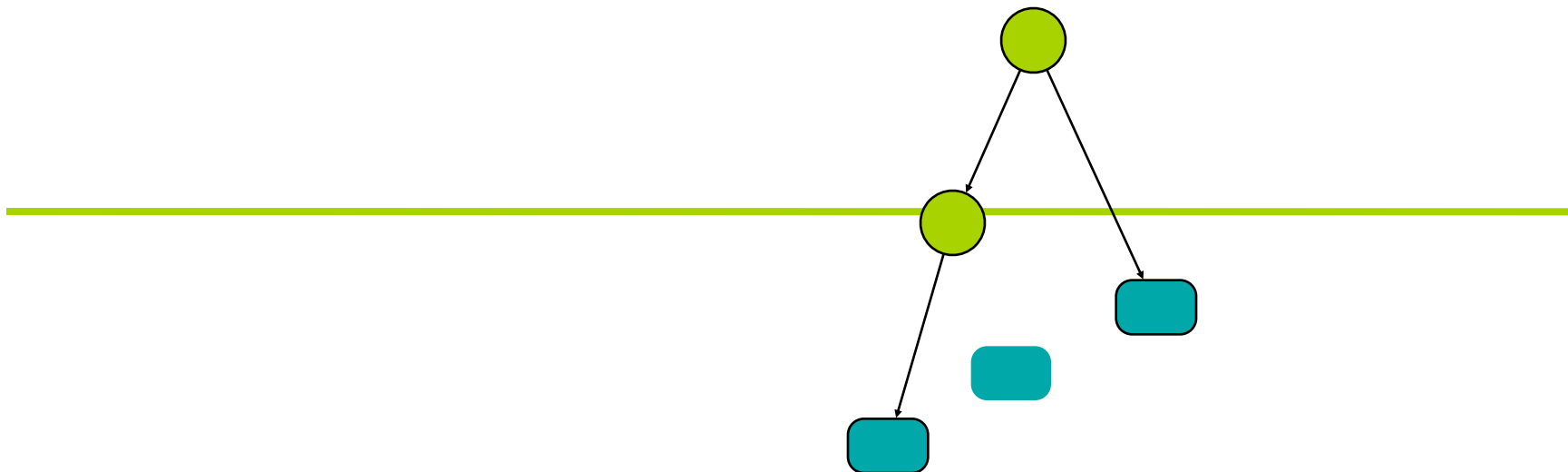
Put the suffix **bab\$** in

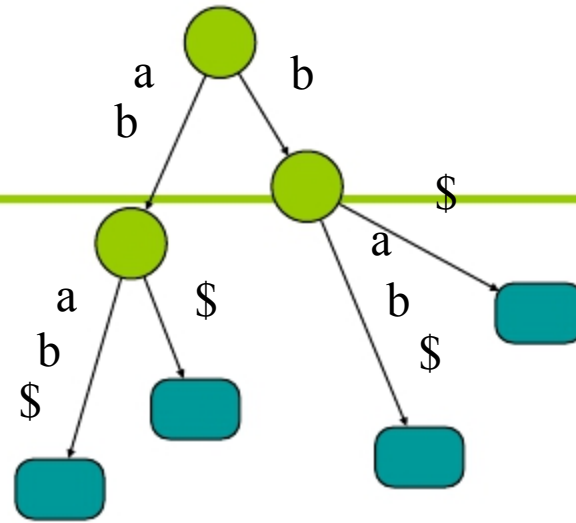




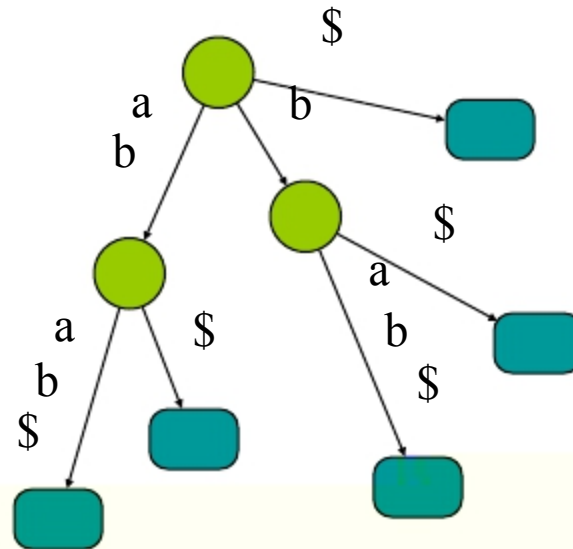
Put the suffi **ab**\$ in

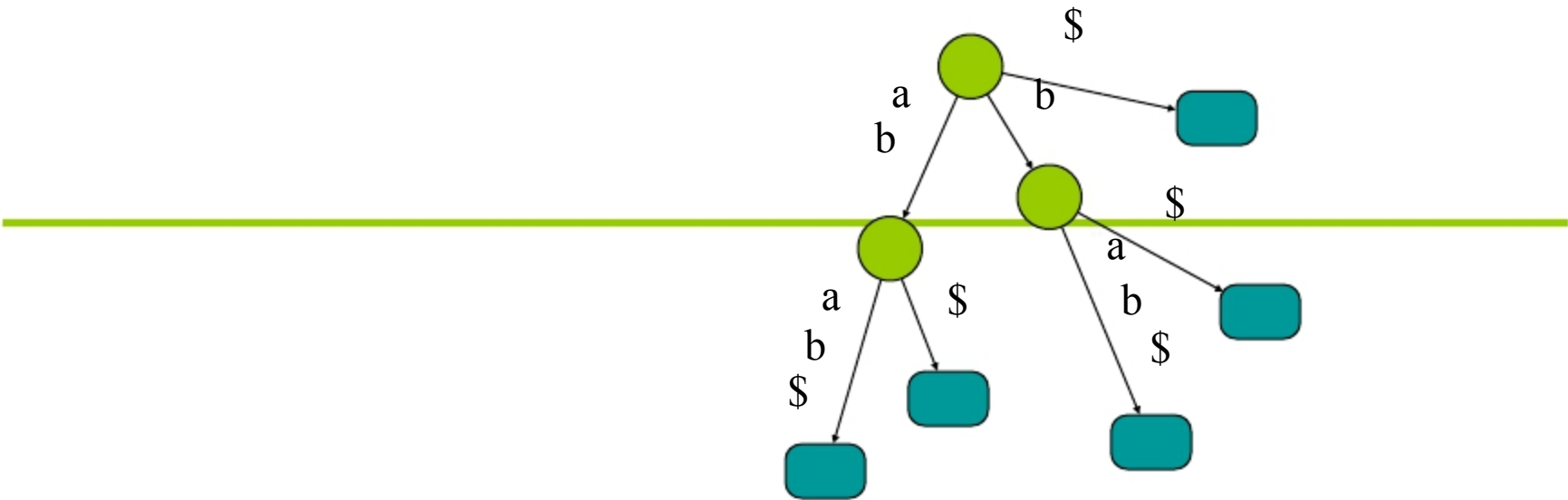




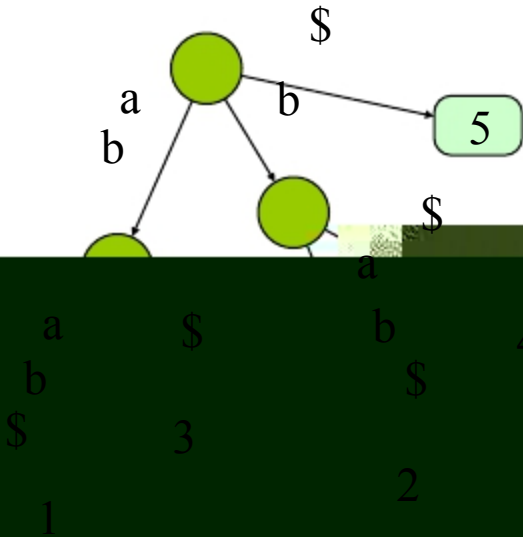


Put the suffi \$ in





We will also label each leaf with the starting point of the corres. suffi .

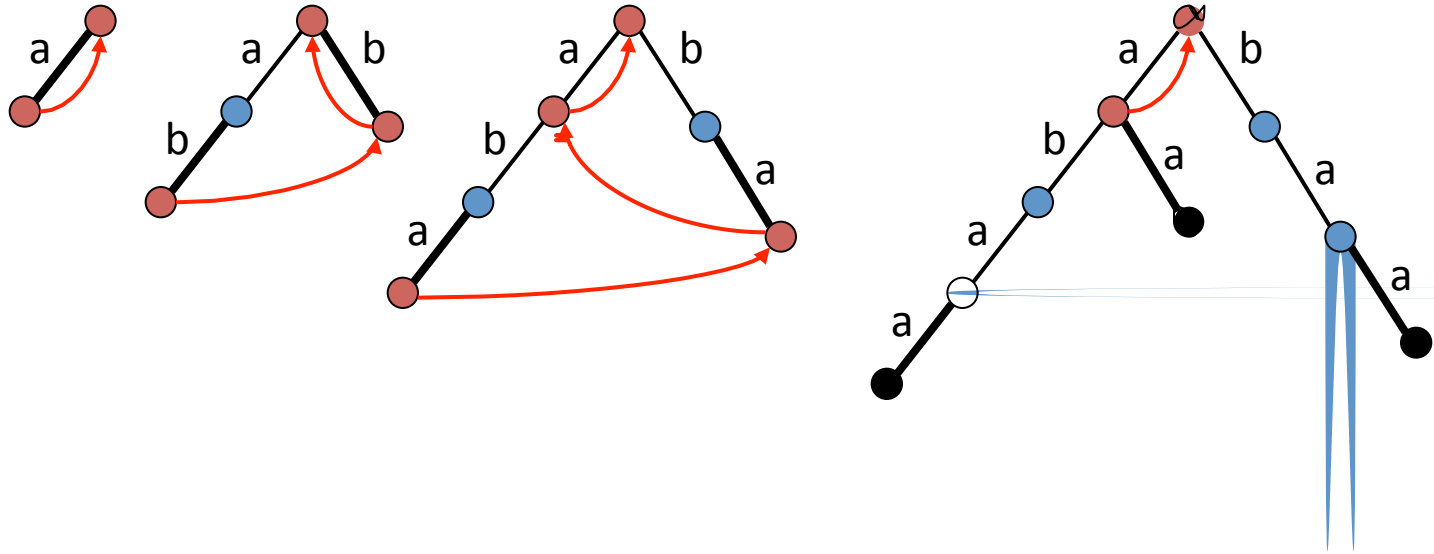


# On-line construction of Trie(T)

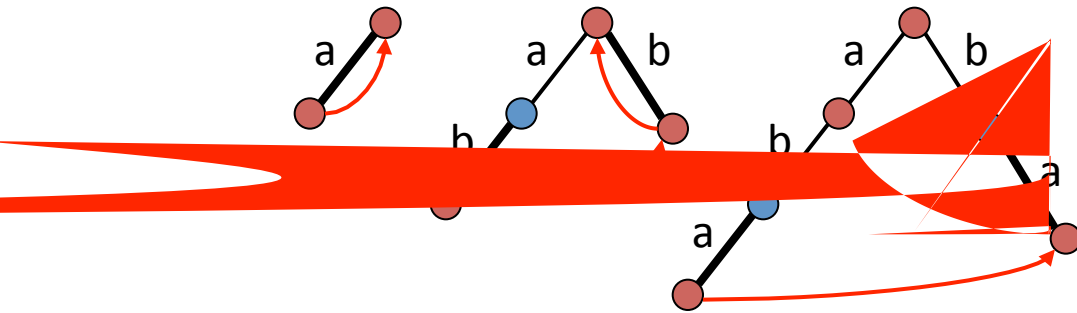
- $T = t_1 t_2 \dots t_n \$$
- $P_i = t_1 t_2 \dots t_i$     *i:th prefix* of  $T$
- on-line idea: update  $Trie(P_i)$  to  $Trie(P_{i+1})$
- $\Rightarrow$  very simple construction



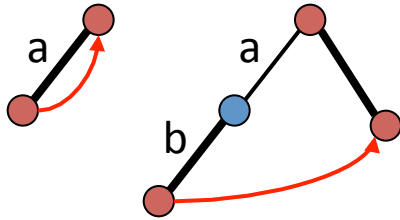
*Trie(abaab)*



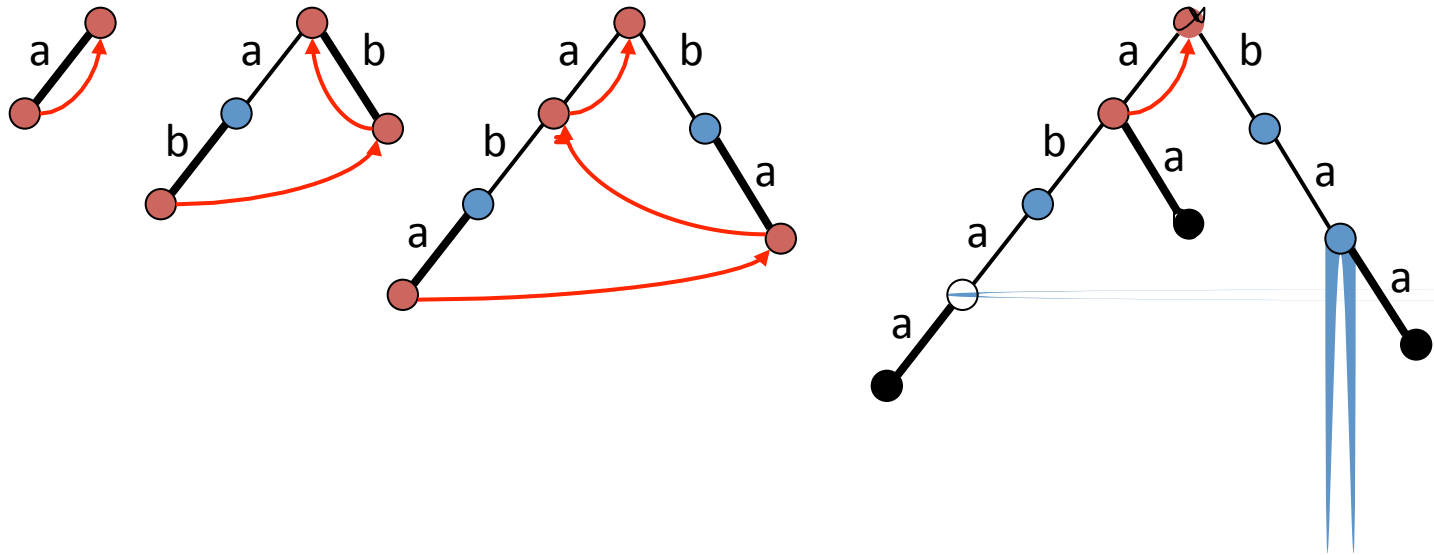
*Trie(abaab)*



# *Trie(abaab)*



*Trie(abaab)*





What happens in  $\text{Trie}(P_i) \Rightarrow \text{Trie}(P_{i+1})$  ?

- time:  $O(\text{size of Trie}(T))$
- suffix links:  
 $\text{slink}(\text{node}(a\alpha)) = \text{node}(\alpha)$

# What can we do with it ?

---

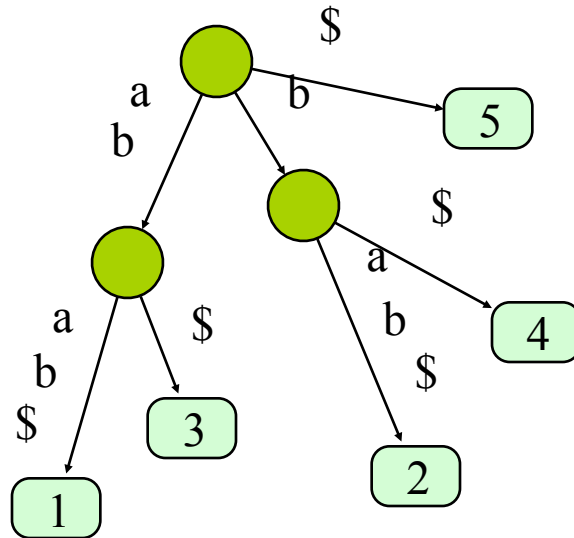
## Exact string matching:

Given a Text  $T$ ,  $|T| = n$ , preprocess it such that when a pattern  $P$ ,  $|P|=m$ , arrives you can quickly decide when it occurs in  $T$ .

We may also want to find all occurrences of  $P$  in  $T$

# Exact string matching

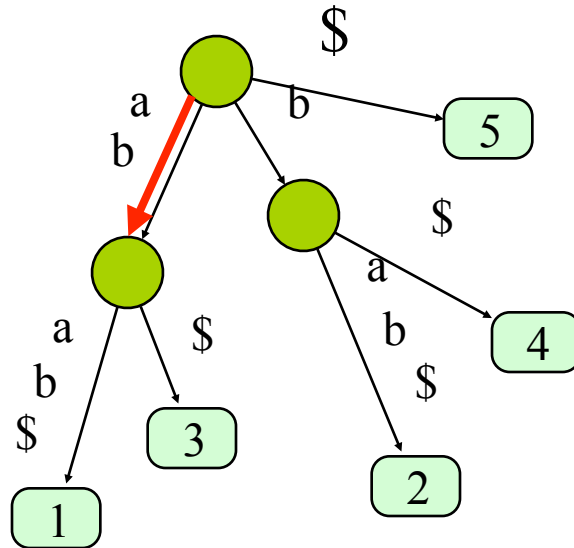
In preprocessing we just build a suffix tree in  $O(n)$  time



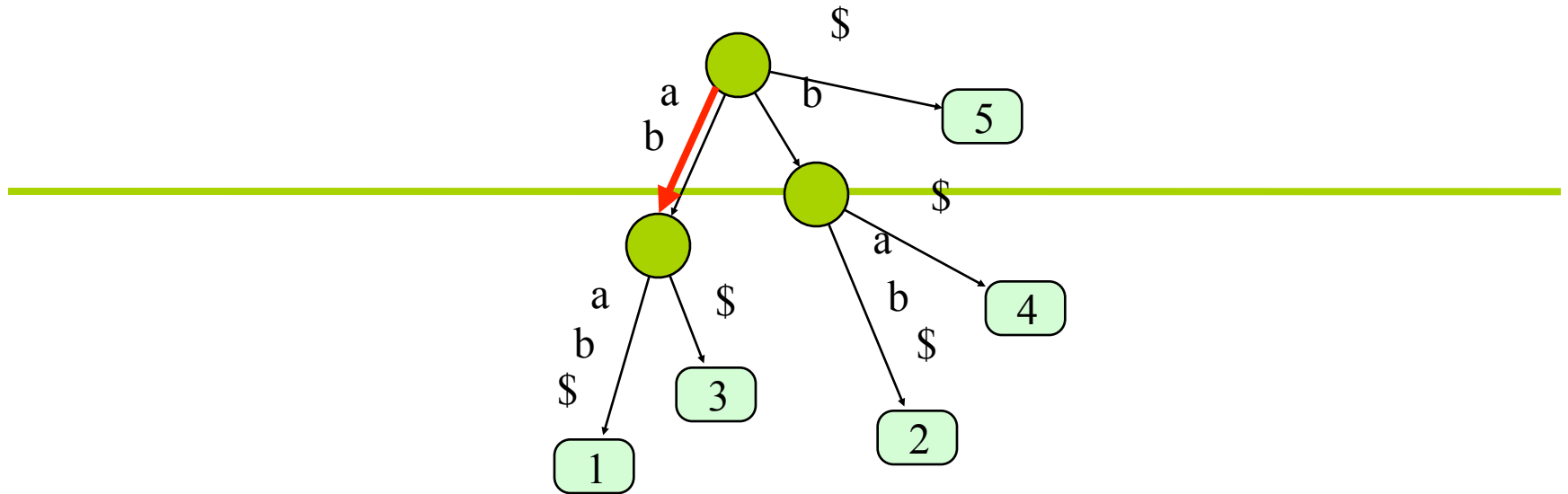
Given a pattern  $P = \text{ab}$  we traverse the tree according to the pattern.

# Exact string matching

In preprocessing we just build a suffix tree in  $O(n)$  time



Given a pattern  $P = \text{ab}$  we traverse the tree according to the pattern.



If we did not get stuck traversing the pattern then the pattern occurs in the text.

Each leaf in the subtree below the node we reach corresponds to an occurrence.

By traversing this subtree we get all  $k$  occurrences in  $O(n+k)$  time

# Generalized suffix tree

---

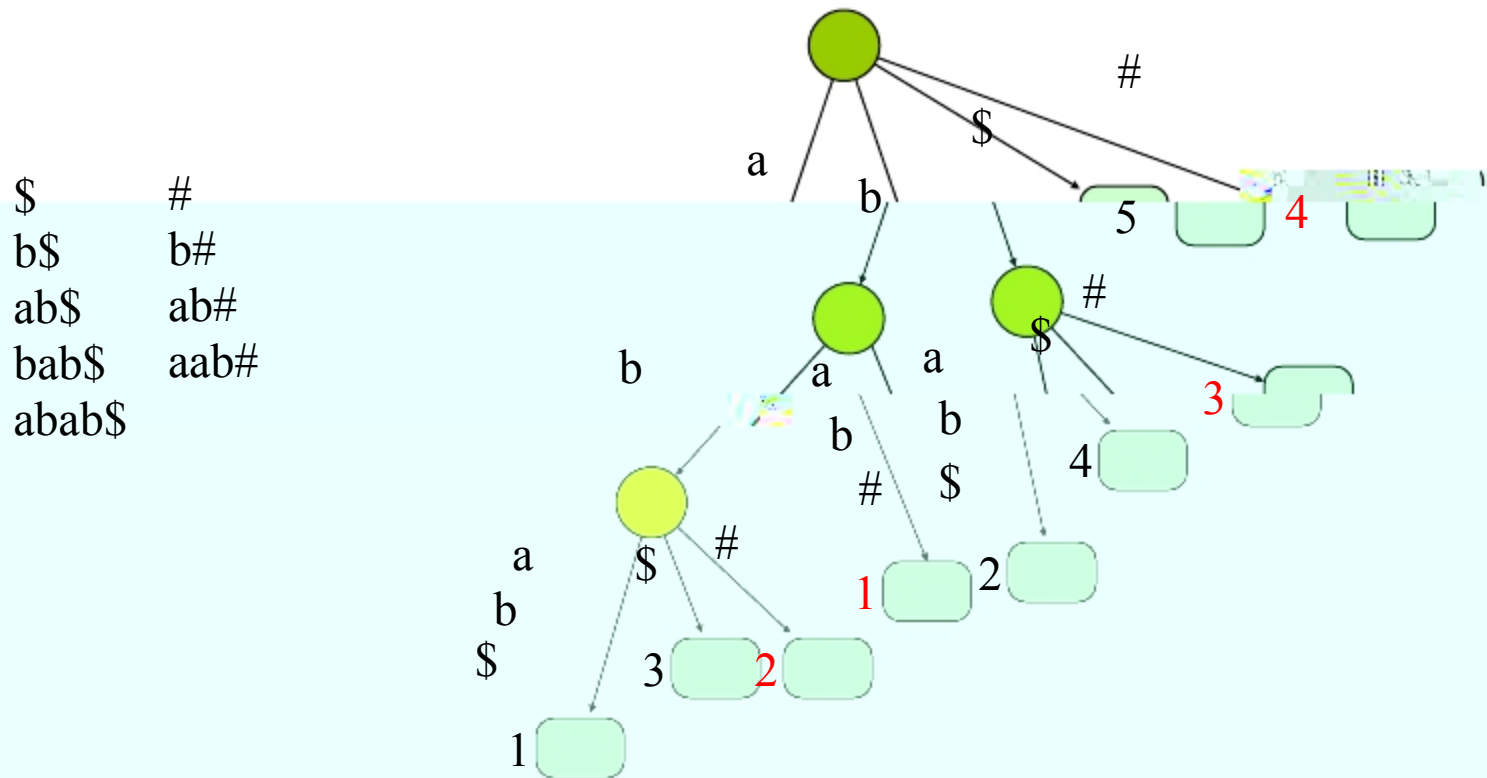
Given a set of strings  $S$  a generalized suffix tree of  $S$  is a compressed trie of all suffixes of  $s \in S$

To make these suffixes prefix-free we add a special char, say  $\$,$  at the end of  $s$

To associate each suffix with a unique string in  $S$  add a different special char to each  $s$

# Generalized suffix tree (Example)

Let  $s_1 = abab$  and  $s_2 = aab$  here is a generalized suffix tree for  $s_1$  and  $s_2$



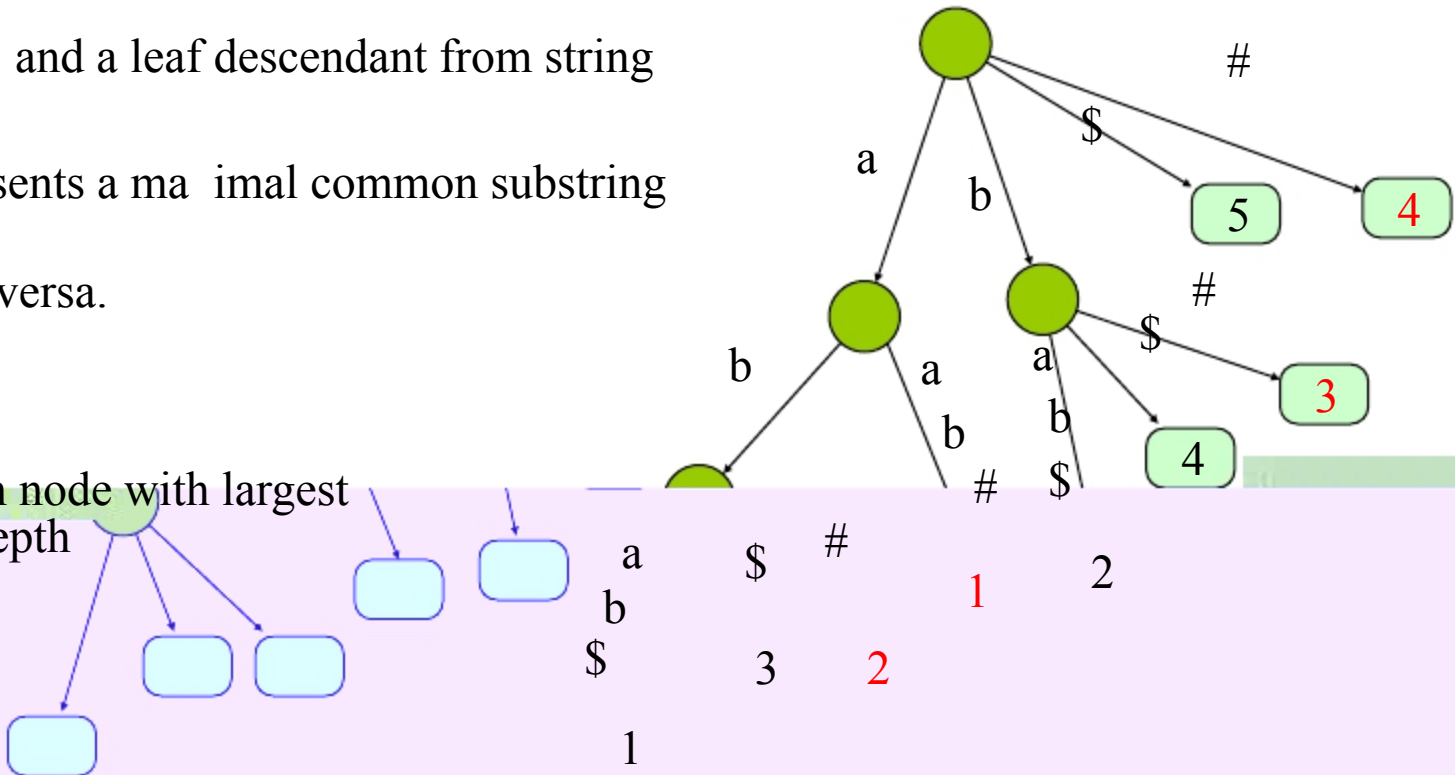
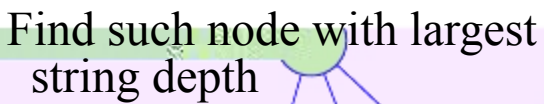
# So what can we do with it ?

---

Matching a pattern against a database of strings

# Longest common substring (of two strings)

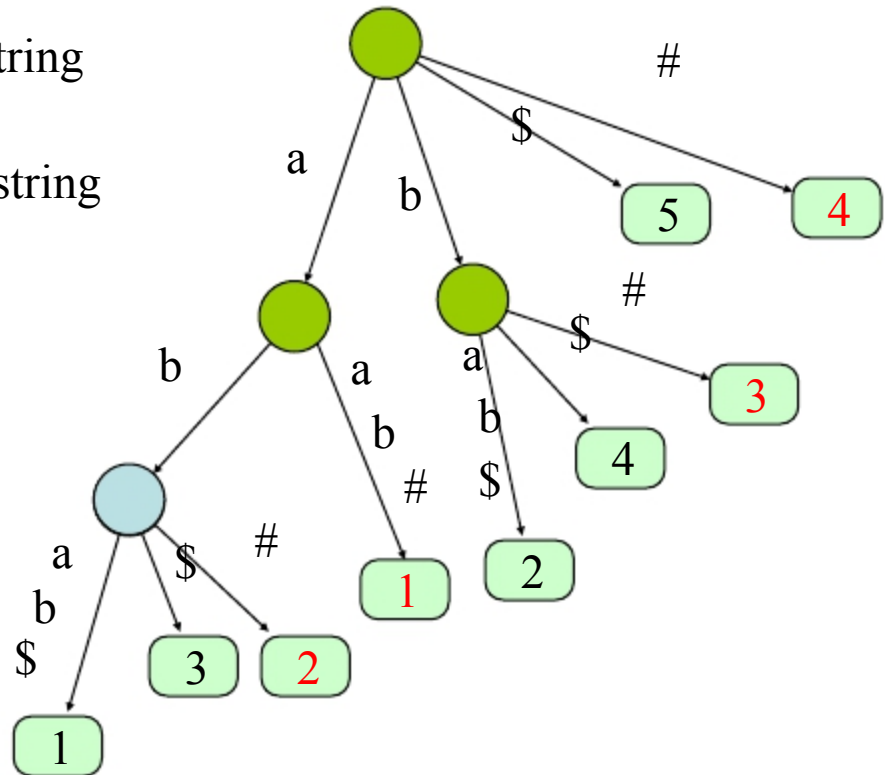
Every node with a leaf descendant from string **S1** and a leaf descendant from string **S2** represents a maximal common substring and vice versa.



# Longest common substring (of two strings)

Every node with a leaf descendant from string **S1** and a leaf descendant from string **S2** represents a maximal common substring and vice versa.

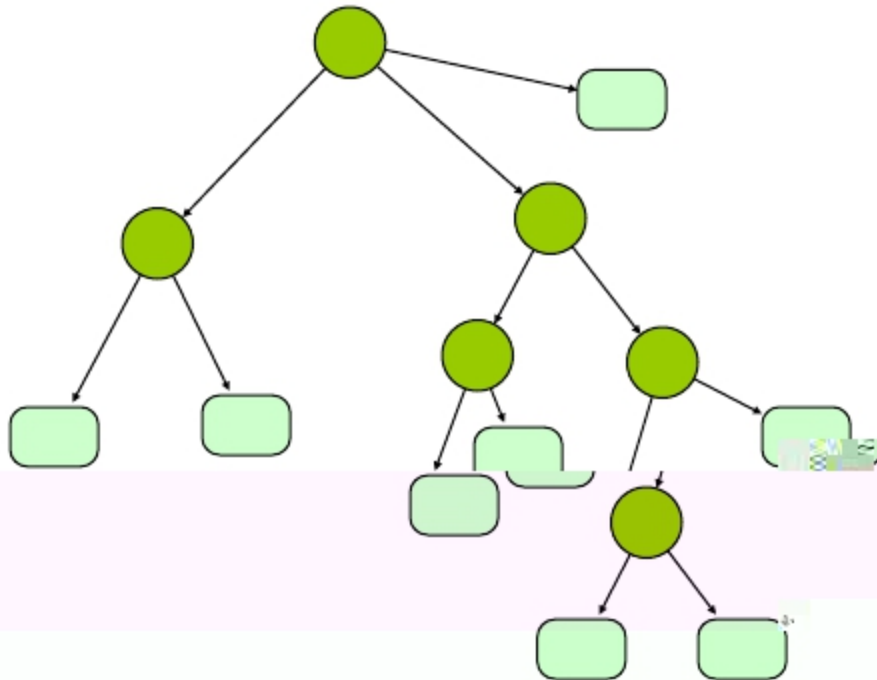
Find such node with largest  
string depth



# Lowest common ancestor

---

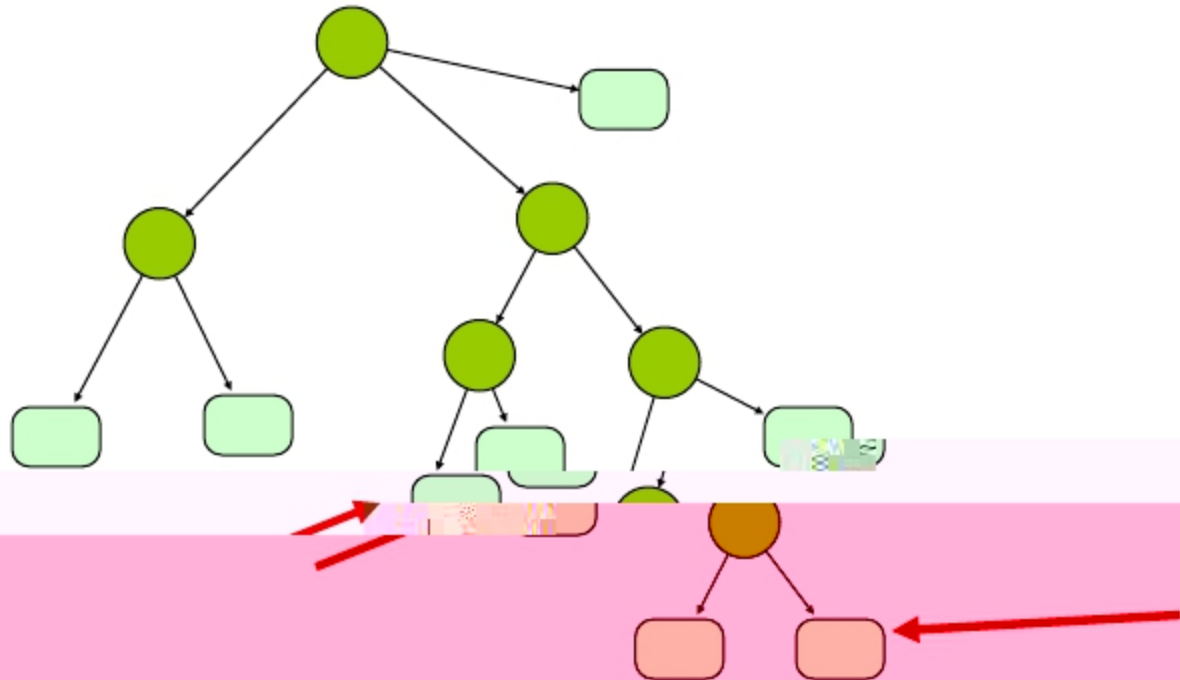
A lot more can be gained from the suffix tree if we preprocess it so that we can answer LCA queries on it



# Lowest common ancestor

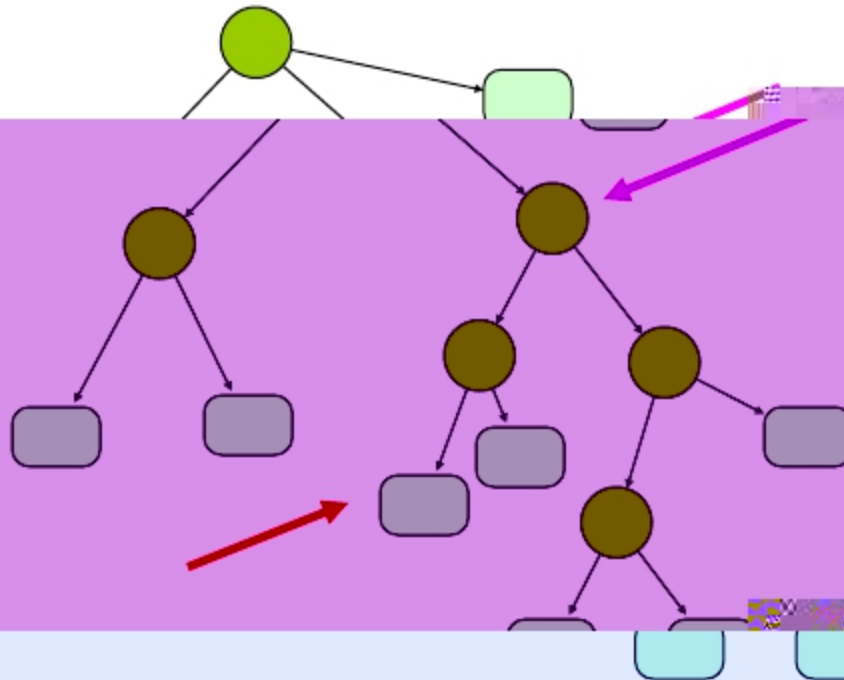
---

A lot more can be gained from the suffix tree if we preprocess it so that we can answer LCA queries on it



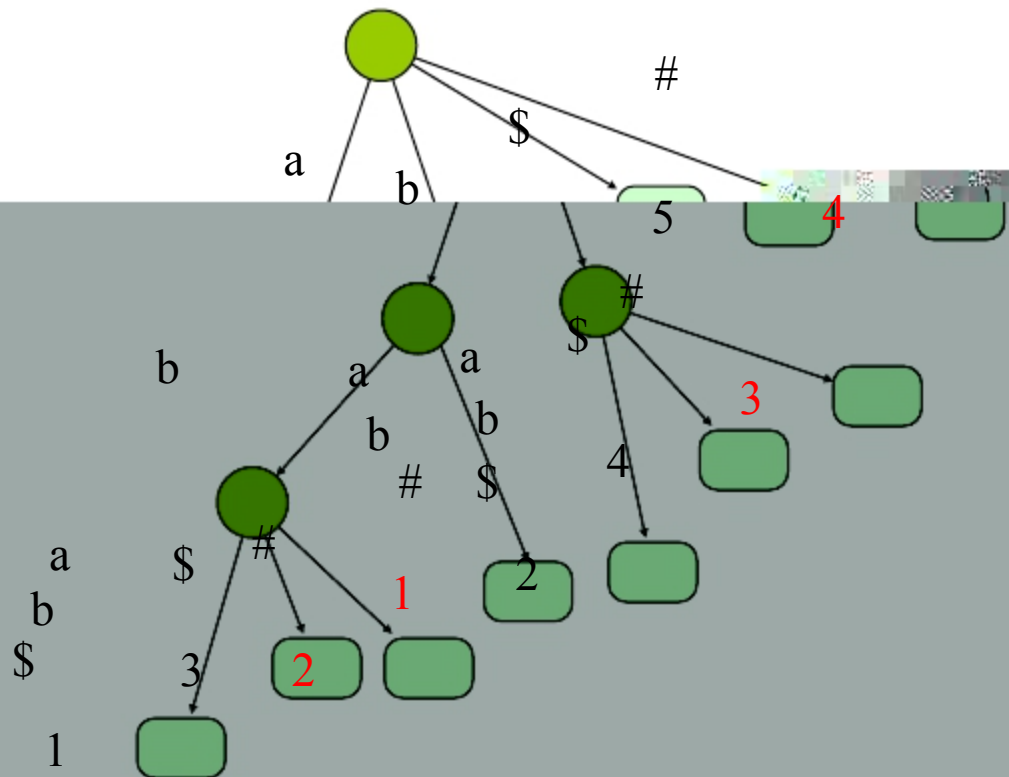
# Lowest common ancestor

A lot more can be gained from the suffix tree if we preprocess it so that we can answer LCA queries on it



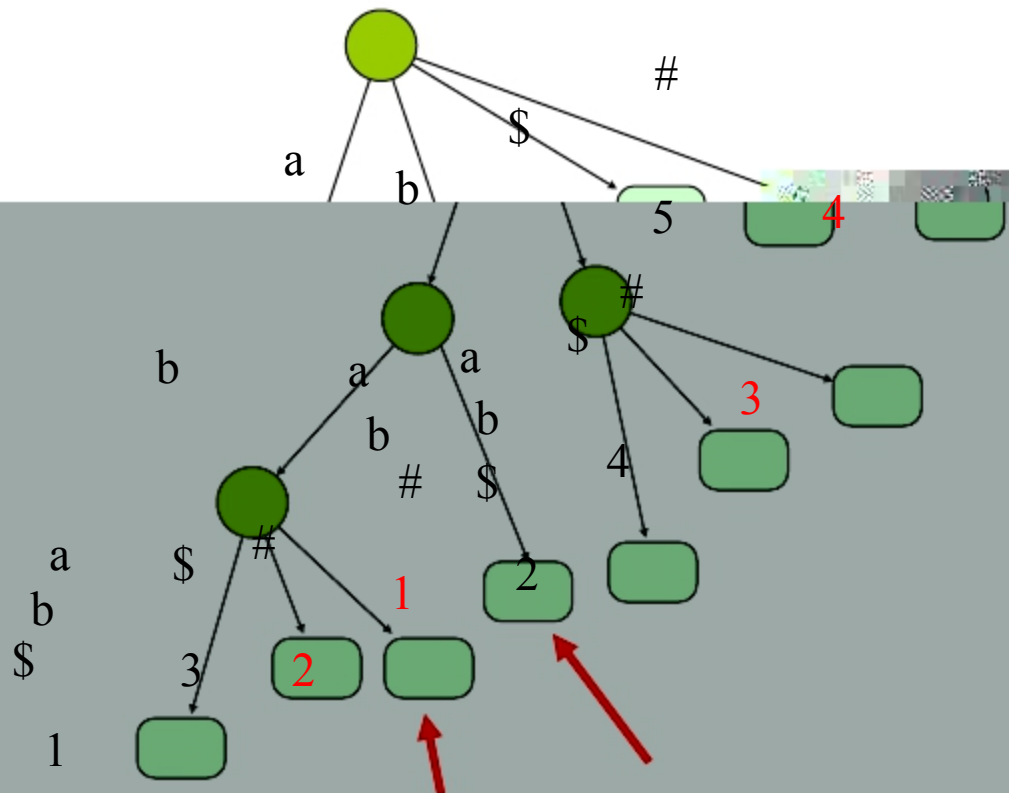
# Wh ?

The LCA of two leaves represents the longest common prefix (LCP) of these 2 suffixes



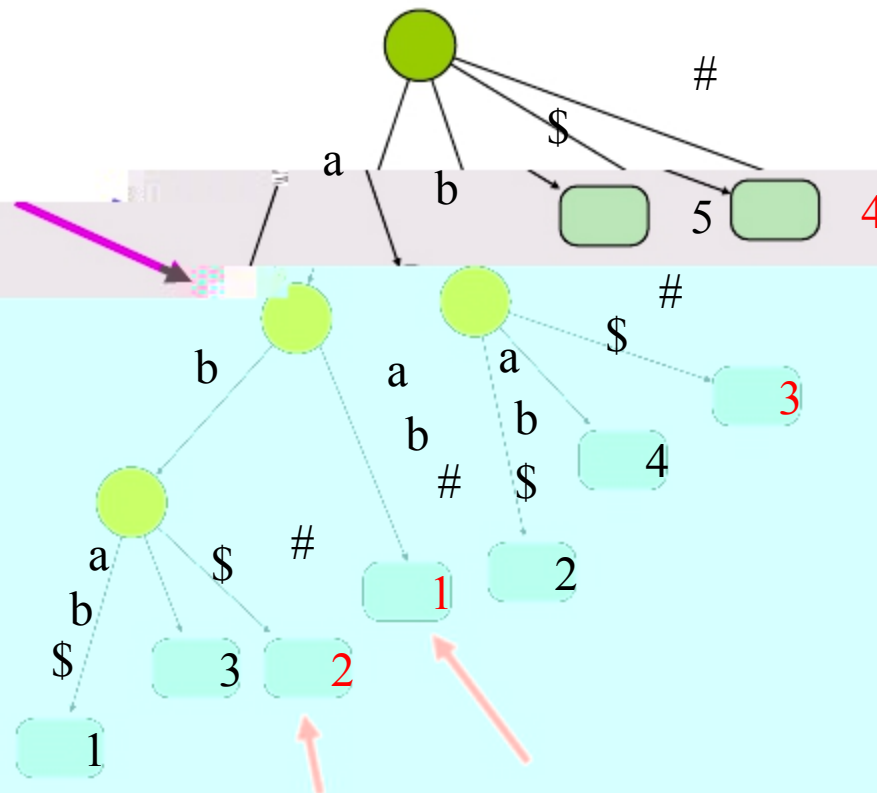
# Wh ?

The LCA of two leaves represents the longest common prefix (LCP) of these 2 suffixes



# Wh ?

The LCA of two leaves represents the longest common prefix (LCP) of these 2 suffixes





# Manimal palindromes algorithm

---

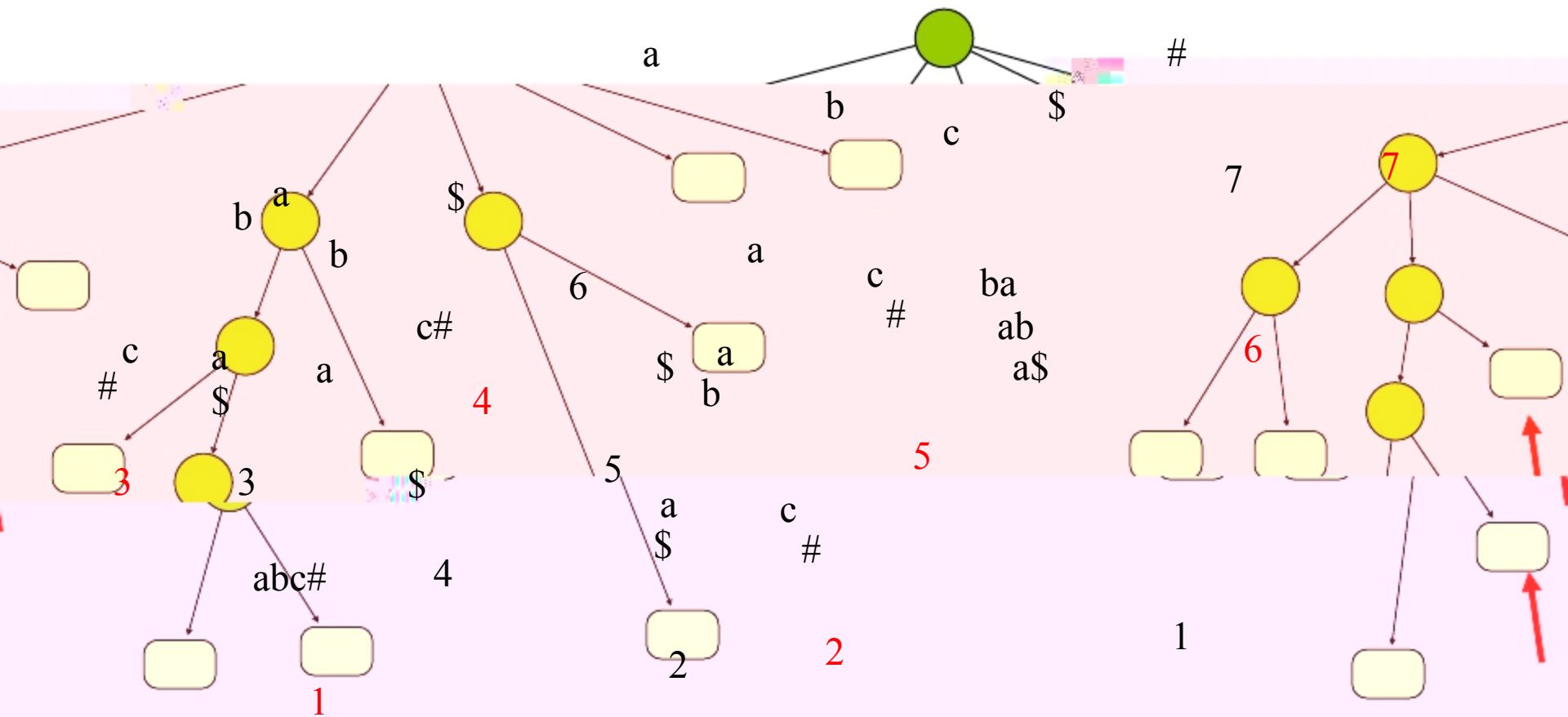
Prepare a generalized suffix tree for

$s = cbaaba\$$  and  $s_r = abaabc\#$

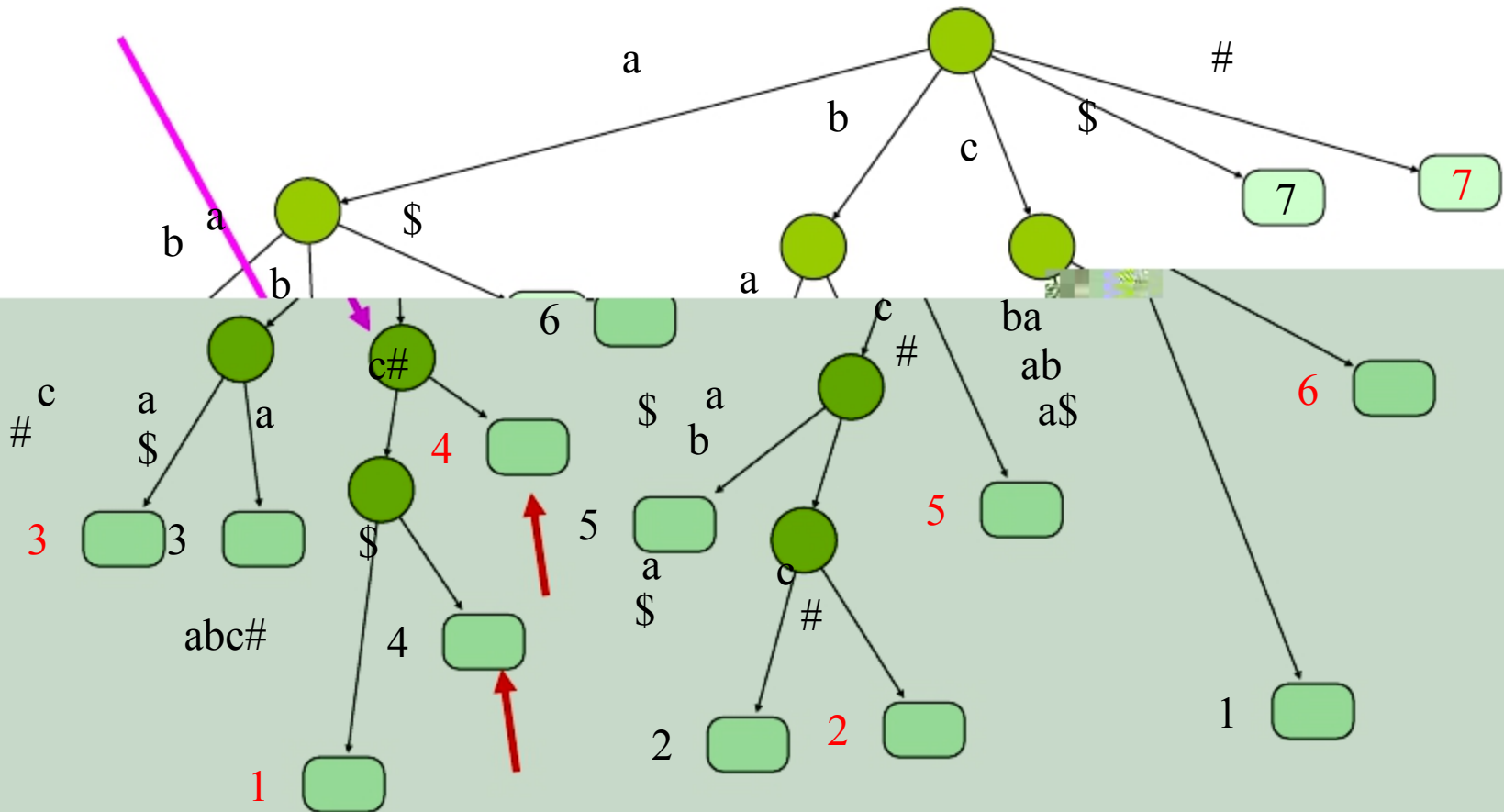
For every  $i$  find the LCA of suffix  $i$  of  $s$  and  
suffix  $m-i+1$  of  $s_r$



Let  $s = cbaaba\$$  then  $s_r = abaabc\#$



Let  $s = cbaaba\$$  then  $s_r = abaabc\#$



# Analysis

---

$O(n)$  time to identify all palindromes

# Drawbacks

---

Suffix trees consume a lot of space

It is  $O(n)$  but the constant is quite big

Notice that if we indeed want to traverse an edge in  $O(1)$  time then we need an array of ptrs. of size  $|\Sigma|$  in each node

# Suffix array

---

We lose some of the functionality but we save space.

Let  $s = abab$

Sort the suffixes lexicographically :  
 $ab, abab, b, bab$

The suffix array gives the indices of the suffixes in sorted order

3	1	4	2
---	---	---	---

# How do we build it ?

---

Build a suffix tree

Traverse the tree in DFS, lexicographically picking edges outgoing from each node and fill the suffix array.

$O(n)$  time

# How do we search for a pattern ?

---

If  $P$  occurs in  $T$  then all its occurrences are consecutive in the suffix array .

Do a binary search on the suffix array

Takes  $O(m \log n)$  time

# E ample

Let **S** = mississippi

---

Let **P** = issa

→	11	i
	8	ippi
	5	issippi
	2	ississippi
	1	mississippi
→	10	pi
	9	ppi
	7	sippi
	4	sisippi
	6	ssippi
→	3	ssissippi

# Supra index

---

Structure

are implementation of

Simple arrays are implemented to the test suffices listed in lexicographical order.

-  
If the suffix array is , this binary search can perform because of the number of random disk accesses.

Suffix arrays are designed to allow done by comparing the contents of each pointer.

To remedy this situation, the use of - over the suffix array has been proposed.

# Supra inde

## E ample

1      6 9 11    17 19   24 28    33      40      46 50    55    60  
This is a te t. A te t has man words. Words are made from letters

<b>60</b>	<b>50</b>	<b>28</b>	<b>19</b>	<b>11</b>	<b>40</b>	<b>33</b>
-----------	-----------	-----------	-----------	-----------	-----------	-----------

Suffi Arra

<b>lett</b>		<b>text</b>		<b>word</b>	
-------------	--	-------------	--	-------------	--

Supra-Inde

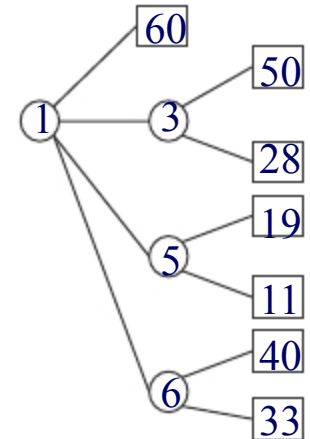
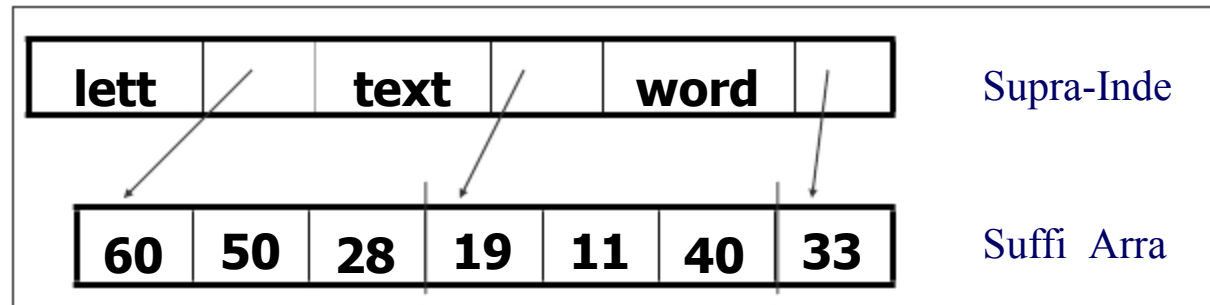
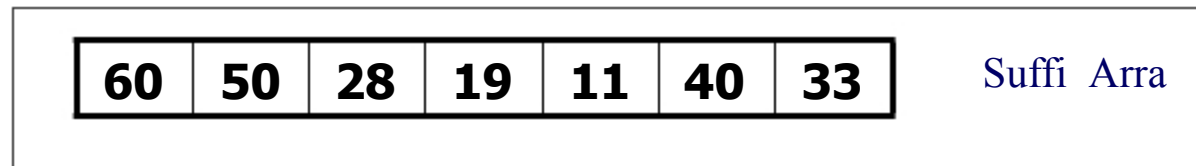
<b>60</b>	<b>50</b>	<b>28</b>	<b>19</b>	<b>11</b>	<b>40</b>	<b>33</b>
-----------	-----------	-----------	-----------	-----------	-----------	-----------

Suffi Arra

# Supra inde

## E ample

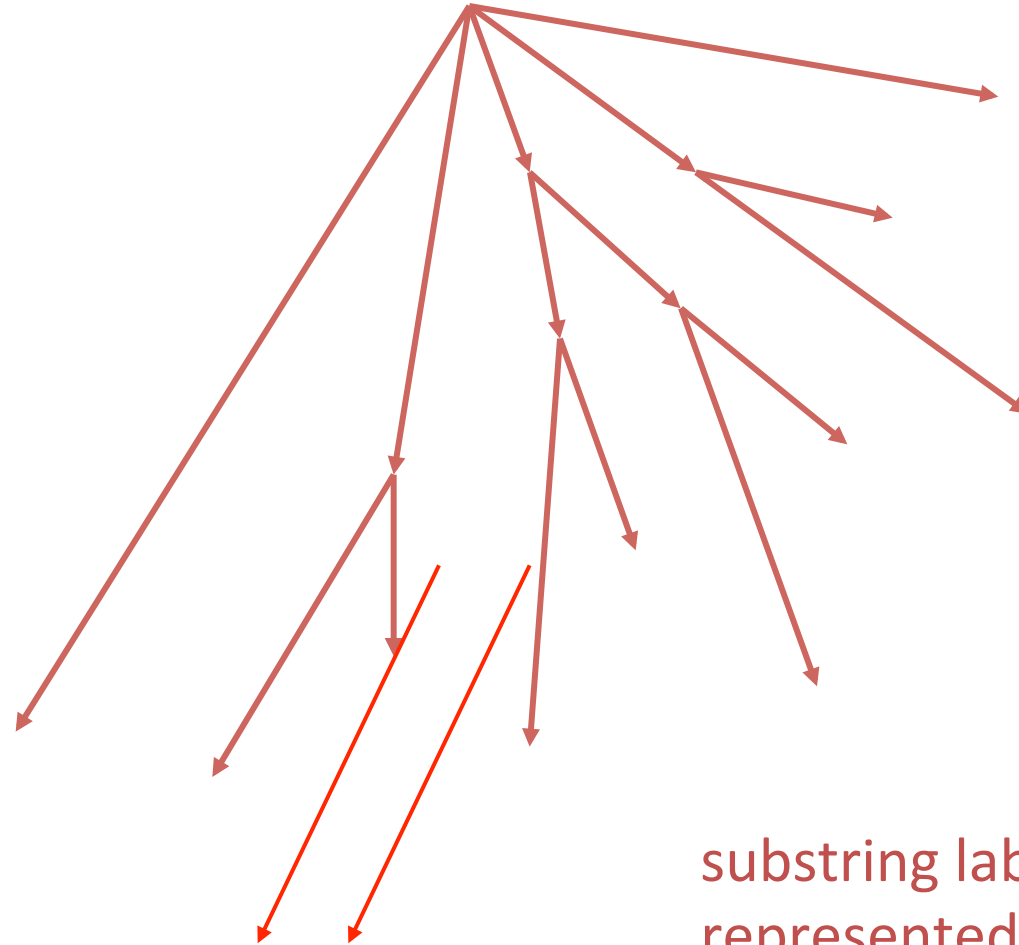
1    6 9 11    17 19    24 28    33    40    46 50    55    60  
 This is a te t. A te t has man words. Words are made from letters



# Tree(hattivatti)

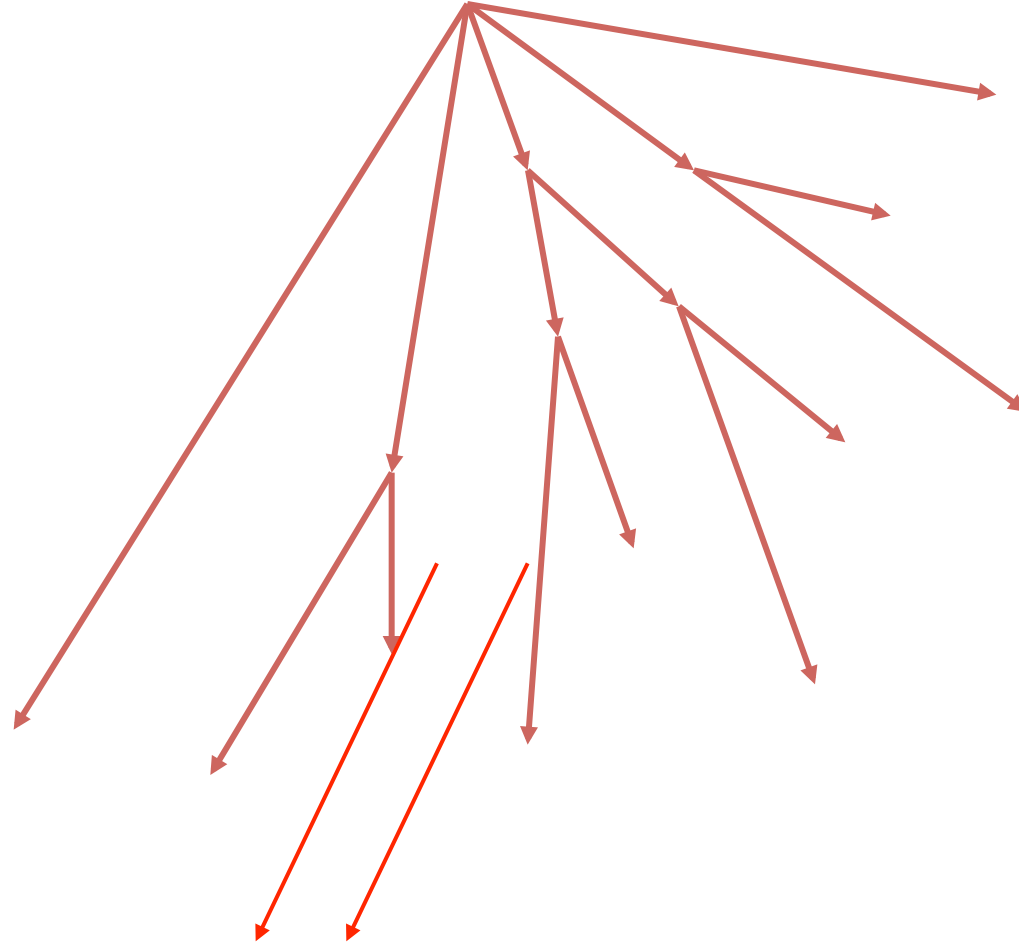


# Tree(hattivatti)

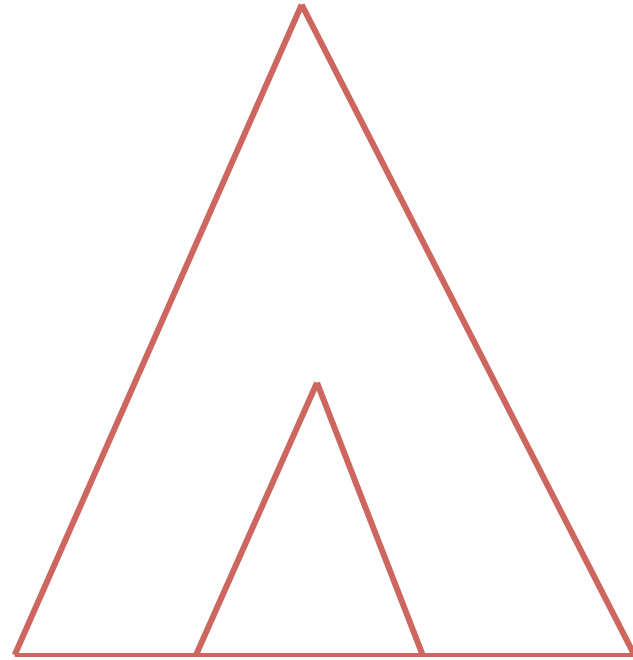


substring labels of edges  
represented as pairs of  
pointers

# Tree(hattivatti)



Tree(T) is *full* text index



Find **att** from Tree(hattivatti)

# Linear time construction of Tree(T)

'on-line' algorithm  
(Ukkonen 1992)



Weiner  
(1973),  
'algorithm  
of the  
year'

McCreight  
(1976)