# 3. Exact Pattern Matching

# String Matching

**Given:** pattern $P$, text $T$    ($P, T$ are strings)

**Aim:** find occurrences of $P$ in $T$

**Exmpl:**

$$T = \overset{\text{1 2 3 4 5 6 7 8 9 10 ...}}{\text{A A T G C A T G C A}} \cdots$$

$$P = \text{A T G}$$

$P$ occurse on position $2, 6, \ldots$

This has applications in:

- Bioinformatics (eg. sequence assembly where one may align fragments of your DNA to reference gnome to get read of your DNA; also applications in finding repeated regions & many more
- general word procemy
- internet search
- "fgrep" in Unix
- search for plagirialism
- subtask for e.g inexact matching

**Basics:**

string $S = x_1 \cdots \cdots x_n$ , $|S| = n$

$$S[i..j] = x_i x_{i+1} \ldots x_j$$

$$S(i) = x_i$$

$S:$ $\overset{\text{1 2 3 4 5 6 7 8}}{\text{H O N O L U L U}}$

$S[1..4]$ prefix

$S[7..6]$ suffix

$S(5) = L$

$S[1..j] \hat{=}$ prefix of $S$ ending at $j$

$S[j..n] \hat{=}$ suffix of $S$ starting at $j$

if $P(i) = T(k)$ for some $i, k$ then they match

else mismatch.

## Naive Method

NAIVE $(P, T)$                // compare $T(i \ldots i + |P| - 1)$
                             // with $P$ for all $i = 1 \ldots |T| - |P| + 1$

    Occurrences $= \emptyset$

    FOR $(i = 1 \ldots |T| - |P| + 1)$        // loop over "pos." in $T$

        FOR $(j = 1 \ldots |P|)$              // loop over $P$

            match $=$ TRUE

            IF $(P(j) \neq T(i + j - 1))$        // compare characters

                match $=$ FALSE

                break

        IF (match)

            add $i$ to occurrences        // save pos. on which
                                                    $P$ occurs

    RETURN occurrences.

Q: How often can P occur in T?
A: |T| - |P| + 1 times

$$T = \overset{1\ 2\ 3\ 4\ 5\ 6}{A\ A\ A\ A\ A\ A}$$
$$P = A\ A$$

P occurs at pos. 1,2,3,4,5
= 5 times = 6 - 2 + 1

Q: What is greatest Nr of character comparisons?
A: |P| · (|T| - |P| + 1)        (worst case)

See example above! 2 · 5 = 10 comparisons

Q: What is least Nr of character comparisons?
A: |T| - |P| + 1        (best case)

$$T = A\ A\ A\ A\ A$$
$$P = B\ A$$

⟹ RUNTIME   NAIVE - alg

$$|T| - |P| + 1 \qquad\qquad 1.\ \text{Loop}$$
$$|P| \qquad\qquad 2\ \text{loop}$$

$$\Longrightarrow \quad |T||P| - |P|^2 + 1 \underset{\underset{|T| \geq |P|}{\uparrow}}{\leq} |T||P|$$

$$\Longrightarrow O(|T||P|)\ \text{time}$$

look now to one of many linear-time
algorithms, i.e. instead of $O(|P| |T|)$

we have runtime $O(|P| + |T|)$

# Z - algorithm [fundamental preprocessing used in many alg.]

**General idea:** pre-process $P$ in $O(|P|)$ time to gain insight of internal structure of $P$

**DEF:** let $S$ be a string (usually $S = P$ pattern) & $i > 1$.

$$Z_i := Z_i(S) = \text{length of longest substring in } S \text{ that starts at position } i \text{ & matches prefix of } S$$

**Example:**

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|-----|---|---|---|---|---|---|---|---|---|----|----|
| $S$ | a | a | b | c | a | a | b | x | a | a | z |
| $Z_i$ | — | 1 | 0 | 0 | 3 | 1 | 0 | 0 | 2 | 1 | 0 |

**DEF (Z-Box):** $\forall\, i > 1$ s.t. $Z_i > 0$, the <u>Z-Box at $i$</u> is the interval $[i, i + Z_i - 1]$

__Exampl:__

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $S$ | a | a | b | c | a | a | b | x | a | a | z |
| $z_i$ | — | 1 | 0 | 0 | 3 | 1 | 0 | 0 | 2 | 1 | 0 |

$z_i > 0$ at pos. $2, 5, 6, 9, 10$

$$z\text{-Box at } 2 = [2, \ 2+1-1] = [2,2]$$
$$\text{at } 5 = [5, \ 5+3-1] = [5,7]$$
$$\text{at } 6 = [6, \ 6+1-1] = [6,6]$$
$$\text{at } 9 = [9, \ 9+2-1] = [9,10]$$
$$\text{at } 10 = [10, \ 10+1-1] = [10,10]$$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $S$ | a | a | b | c | a | a | b | x | a | a | z |

this intervals $[i,j]$ correspond
to $S[i..j] =$ longst substring
starting at $i > 1$ &
matches prefix of $S$

<u>DEF</u>   $\forall i > 1$ :

(r = right)   $r_i$  denotes right-most endpoint of
z-Boxes $z_j$ with $z_j > 0$ & $j \leq i$

(equ. to $r_i$ = largest value of $j + z_j - 1$
over all $1 < j \leq i$ st $z_j > 0$)

"store index j": $l_i = j$  for  $j$  satisfying this ↗

(l = left)   (equ. to $l_i$ = position of left-end
of z-Box ending in $r_i$ )

if MORE than one z-Box ends in $r_i$
then $l_i$ can be chosen arbitrarily among those values

<u>Formal:</u>   $r_i = \max\limits_{2 \leq j \leq i} \{ j + z_j - 1 : z_j > 0 \}$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $s$ | a | a | b | c | a | a | b | x | a | a | z |
| $r_i$ | – | 2 | 2 | 2 | 7 | 7 | 7 | 7 | 10 | 10 | 10 |
| $l_i$ | – | 2 | 2 | 2 | 5 | 5 | 5 | 5 | 9 | 9 | 9 |

right:   $r_6 = \max \{ \underbrace{2+1-1}_{\text{for } j=2}, \underbrace{5+3-1}_{j=5}, \underbrace{6+1-1}_{j=6} \} = 7$

left:   $l_6 = 5$

this essentially gives pos. of "large" z-Boxes.

## How to compute $z_i, r_i, l_i$ efficiently?

_main idea:_ iterative approach

start with $z_2$ (explicit comparison from left-to-right)

Assume $\forall i < k$ values $z_i, r_i, l_i$ have been computed.

for $k$ use: $z_i, l_i, r_i$ $(i < k)$
WITHOUT EXPLICITE character comparisons
as much as possible.

## Overview of steps of Z-alg

①    compute $z_2$ (explizite comparison of $S[1...n]$ with $S[2...n]$ until mismatch
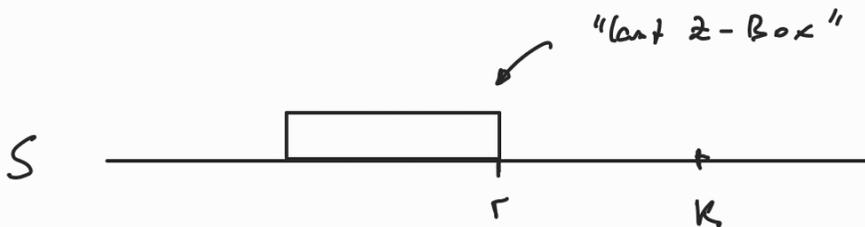
IF $(z_2 > 0)$ put $l_2 = 2$, $r_2 = z_2 + 1$
ELSE           put $l_2 = r_2 = 0$

(2)  $k$-th step  $(k \geq 2)$.

$\Rightarrow$  all  $z_i, l_i, r_i$  computed for all $i \leq k-1$

Compute $z_k, l_k, r_k$  based on the
following cases.

$$l := l_{k-1}, \quad r := r_{k-1}$$

Case 1:    $r < k$



"last $z$-Box"

$S$

$r$      $K$

in this case  compute $z_k$ via
explicit comparison  of  $S[k..n]$ & $S[1...n]$
                until mismatch.

IF $(z_k > 0)$ put  $l_k = k$,  $r_k = k + z_k - 1$
ELSE                    $l_k = l$,  $r_k = r$

**Exmpl:**

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $S$ | a | a | b | c | a | a | b | x | a | a | z |
| $r_i$ | – | 2 | 2 | 2 | 7 | 7 | 7 | 7 | | | |
| $l_i$ | – | 2 | 2 | 2 | 5 | 5 | 5 | 5 | | | |

$$k = 9 \longrightarrow r_8 = 8 < k = 9 \quad \text{ned to compare} \ S[9..11] \ \text{with} \ S[1..n]$$

$$\Rightarrow \quad S[9] = S[1] \ \checkmark$$
$$S[10] = S[2] \ \checkmark$$
$$S[11] \neq S[3] \ X$$

$$\Rightarrow \quad z_9 = 2, \quad l_9 = 9$$
$$r_9 = 9 + 2 - 1 = 10$$

## Case 2: $r \geq k$

$$S \quad\rule{6cm}{0.4pt}$$

z-Box over $[\ell, k, r]$, with $z_k$ to the right.

length $r - \ell + 1 = z_\ell$

z-Box $\hat{=}$ prefix of $S$

$\Longrightarrow$

same!

$S$

positions: $1 \quad k' \quad z_\ell \qquad \ell \quad k \quad r$

with boxes $\alpha$, $\beta$ repeated.

$$S[1 .. z_\ell] \quad = \quad S[\ell .. r] = \alpha$$

$$S[k' .. z_\ell] \quad = \quad S[k .. r] = \beta$$

We want to know now $z_k$ ( longest string starting at $k$ $\alpha$ is prefix of $S$ )

Question: what we know about $\beta$ ?

if we know $\beta$ or "first part" of $\beta$ is prefix of $S$, we don't need to compare the respective positions.

! This knowledge is already stored in $z_{k'}$

$$\boxed{k' = k - \ell + 1}$$

and leads to following subcases.

$$2A, \quad 2B$$

$$\boxed{2A} \quad z_{k'} < |\beta| = r - k + 1$$



identically

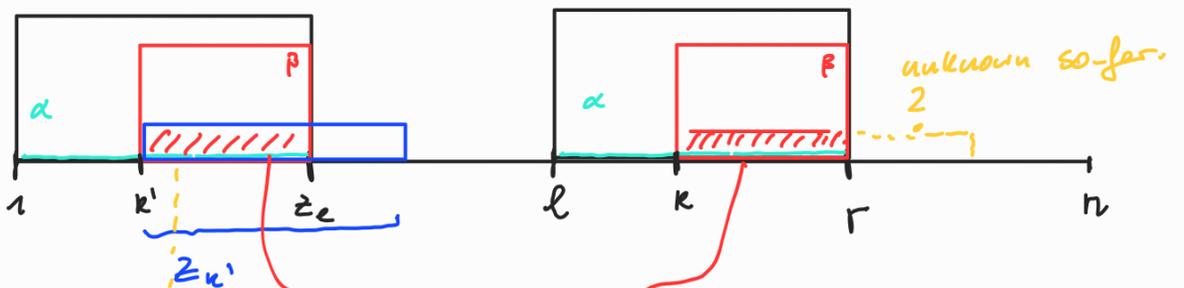$z_{k'} =$ length of $\gamma$ = length of longest substring starting at $k'$ & is prefix of $S$

Since $\beta = \beta \atop \text{left} \quad \text{right}$ $\Rightarrow$ $\gamma_{(right)}$ is longest substring starting at $k$ & is prefix of $S$.

$$\Rightarrow \quad z_k = z_{k'}$$

$l$ & $r$ remain unchanged.

$$\boxed{2B} \quad z_{k'} \geq |\beta| = r - k + 1$$



unknown so-far.
?

$z_{k'}$

$\beta + 1$
pos.

this part matches already, id. $\beta_{(right)}$ is part of longest substring starting at $k$ & prefix of $S$.

$\beta = S[k..r]$ is prefix of $S$.

& $z_k \geq |\beta| = r - k + 1$

However $Z_k > |\beta|$ may be possible.

$\Rightarrow$ compare $S[r+1..n]$ [yellow part above]

with $S[|\beta|+1,...n]$ until mismatch

$r - k + 1 + 1$
$= r - k + 2$

to possibly extend
Z-Box starting at $k$

let $q$ be pos. in $S$ of 1st mismatch
& if no mismatch put $q = |S| + 1$

put $Z_k = q - k$ $(= k + Z_k - 1)$
$r = q - 1$
$\ell = k$

## Z Algorithm

1: $r \leftarrow \ell \leftarrow 0$
2: **for** $k = 2$ to $|S|$ **do**
3: // Case 1:
4:    **if** $r < k$ **then**
5:       Compare characters in $S[k..n]$ with the ones in $S[1..n]$ until mismatch is found (from left-to-right).
6:       Set $Z_k$ to the length of the matched characters
7:       **if** $Z_k > 0$ **then**
8:          Set $r \leftarrow k + Z_k - 1$, $\ell \leftarrow k$.
9: // Case 2:
10:    **else**
11:       $k' \leftarrow k - \ell + 1$
12: // Case 2a:
13:       **if** $Z_{k'} < r - k + 1$ **then**
14:          $Z_k \leftarrow Z_{k'}$
15: // Case 2b:
16:       **else**
17:          $\ell \leftarrow k$
18:          Compare characters in $S[r + 1..n]$ with the ones in $S[r - k + 2..n]$ until mismatch is found (from left-to-right).
19:          let $q > r$ be the position of the first mismatch or $q = |S| + 1$ if no mismatch occurs
20:          $Z_k \leftarrow q - k$, $r \leftarrow q - 1$

latter discussion implies:

Thm: $z$-alg. correctly computes all
$Z$-Box $z_i$, $i > 2$

However more important:

Thm: All $z_i(S)$ values are computed in $O(|S|)$ time.

proof:

For loop (Line 2-20) runs $O(|S|)$ times.

Q: What happens within FOR-loop?

to answer this question let us count
character-comparisons.

Each character comparison results either in match
or mismatch.

$\Rightarrow$ let us count match / mismatches.

each comparison ends when 1st mismatch occurs
$\Rightarrow$ since $O(|S|)$ different comparisons
$\Rightarrow$ total $O(|S|)$ mismatches

Notation: $c_k =$ NR of comparisons in $k$-th iteration

$m_k =$ NR of matches in $k$-th iteration

Claim: $r_k - r_{k-1} \geq m_k$

Proof: in $z$-Alg we either are Case 1/2A/2B.

Assume CASE 1 applies:

& thus, $r_{k-1} < k$ $\xrightarrow{\text{Alg.}}$ explicit comparison of $S[k..n]$ with $S[1..m]$ until mismatch.

$\Rightarrow$ at most $n-k+1$ comparisons

in Alg we put $z_k = m_k$

in Alg case $z_k > 0$ & [$z_k = 0$ implicit by leaving $"r = r_k = r_{k-1}"$ unchanged].

$z_k > 0$: $r_k = k + z_k - 1 = k + m_k - 1 \geq r_{k-1} + m_k$

$\underset{m_k}{"}$

$k-1 \geq r_{k-1}$

$\iff r_k \geq r_{k-1} + m_k$

$\implies r_k - r_{k-1} \geq m_k$ ✓

$z_k = 0$: $r_k = r_{k-1} = r_{k-1} + \underset{0}{\overset{||}{m_k}}$

(implicit in algo)

$\implies r_k - r_{k-1} = m_k = 0.$ ✓
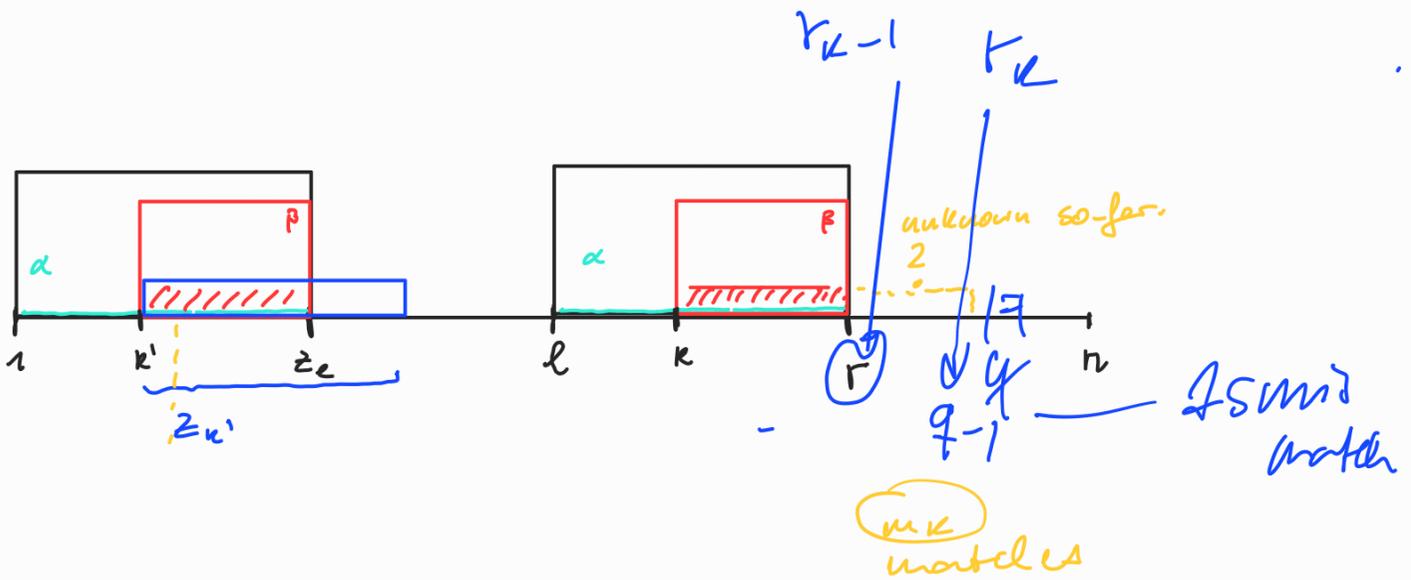
Now Case 2A/2B.

CASE 2A:   in Alg $\overset{=}{r}=r_{k-1}$ remains unchanged,

i.e   $r_k = r_{k-1}$

moreover   no comparisons are
made,   i.e.,   $m_k = 0$

$\implies$   $r_k - r_{k-1} = m_k$   ✓

CASE 2B:   in Alg :  $q > r_{k-1}$ & $r_k = q - 1$ :



$\implies$   $r_k - r_{k-1} \geq m_k$   [see picture] . ✓

$n := |S|$

$\implies$   $\displaystyle\sum_{k=2}^{n} c_k \;\leq\; \sum_{\substack{k=2 \\ (mism.)}}^{n} 1 \;+\; \sum_{\substack{k=2 \\ (match)}}^{n} m_k \;=\; n-1 + \sum_{k=2}^{n} m_k$

$\leq n-1 + (r_2 - r_1) + (r_3 - r_2) + (r_4 - r_3) \cdots + (r_n - r_{n-1}) = n - 1 + \underbrace{r_n}_{\leq n} - \underbrace{r_1}_{=0} \leq 2n-1 \in O(n)$   □

(claim)

in Algo $r_1 = \ell_1 = 0$

## Exmpl

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-----|---|---|---|---|---|---|---|---|---|----|
| $S$ | a | c | a | c | a | b | a | c | a | c |
| $z_i$ | – | | | | | | | | | |
| $l_i$ | – | | | | | | | | | |
| $r_i$ | – | | | | | | | | | |

Init $r = l = 0$

$k = 2$ : $\quad k = 2 > r = 0 \quad$ ( Case 1)

$\qquad$ Compare $S[2\cdots n]$ with $S[1..n]$

$\qquad\qquad \Rightarrow S(2) \neq S(1)$ mismatch

$\qquad\qquad \Rightarrow z_k = 0 \quad,\quad r, l$ unchanged.

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-----|---|---|---|---|---|---|---|---|---|----|
| $S$ | a | c | a | c | a | b | a | c | a | c |
| $z_i$ | – | 0 | | | | | | | | |
| $l_i$ | – | 0 | | | | | | | | |
| $r_i$ | – | 0 | | | | | | | | |

$k = 3.$ , $k = 3 > r = 0 \quad$ ( Case 1)

$\qquad$ Compare $S[3\cdots n]$ with $S[1\cdots n]$

$\qquad\qquad S(3) = S(1)$
$\qquad\qquad S(4) = S(2)$
$\qquad\qquad S(5) = S(3)$
$\qquad\qquad S(6) \neq S(4)$

$\Rightarrow$

$\boxed{\begin{aligned} z_k &= 3 > 0 \\ r &= 3 + 3 - 1 = 5 \\ l &= 3 \end{aligned}}$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $S$ | a | c | a | c | a | b | a | c | a | c |
| $z_i$ | – | 0 | 3 | | | | | | | |
| $\ell_i$ | – | 0 | 3 | | | | | | | |
| $r_i$ | – | 0 | 5 | | | | | | | |

$$k=4 \quad, \quad r=5 \geqslant k=4 \quad (\text{case } 2)$$

$$(\alpha = aca, \quad \beta = c\alpha)$$

$$k' = k - \ell + 1 = 4 - 3 + 1 = 2$$

$$z_{k'} = z_2 = 0 < r - k + 1 = 5 - 4 + 1 = 0$$

$$\Rightarrow \text{ case } 2A: \quad z_4 = z_{k'} = 0 \quad, \quad r, \ell \text{ unchanged.}$$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $S$ | a | c | a | c | a | b | a | c | a | c |
| $z_i$ | – | 0 | 3 | 0 | | | | | | |
| $\ell_i$ | – | 0 | 3 | 3 | | | | | | |
| $r_i$ | – | 0 | 5 | 5 | | | | | | |

$$k=5 \quad, \quad r=5 \geqslant k \quad (\text{case } 2 \text{ with } \alpha = aca, \ \beta = a)$$

$$k' = k - \ell + 1 = 5 - 3 + 1 = 3$$

$$z_{k'} = z_3 = 3 > r - k + 1 = 5 - 5 + 1 = 1$$

$$\Rightarrow \text{ case } 2A \text{ again and so on}\cdots$$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-----|---|---|---|---|---|---|---|---|---|----|
| $S$ | a | c | a | c | a | b | a | c | a | c |
| $z_i$ | – | 0 | 3 | 0 | 1 | 0 | 4 | 0 | 2 | 0 |
| $l_i$ | – | 0 | 3 | 3 | 5 | 5 | 7 | 7 | 9 | 9 |
| $r_i$ | – | 0 | 5 | 5 | 5 | 5 | 10 | 10 | 10 | 10 |

Based on preprocessing with $z$-alg. we can design a:

## Simple linear-time exact matching Alg

Simple-exact-matching $(P, T, \$)$    // $\$$ character not in $T$ & $P$

1   occurrences = $\phi$
2   $n = |P|$,   $m = |T|$
3   $S = P \$ T$
4   preprocess $S$ with $Z$-Alg to compute $z_2 \dots z_{|S|}$
5   FOR $(i = n+2, \dots, |S| = m+n+1)$
6     IF $(z_i = n)$
7      add $i-n-1$ to occurrences

_Exmpl_     $P = bbac$  ,   $T = abba\ bbaca$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $S$ | b | b | a | c | $\$$ | a | b | b | a | b | b | a | c | a |
| $z_i$ | (doesn matter..) | | | | | 0 | 3 | 1 | 0 | 4 | 1 | 0 | 0 | 0 |

$z_{10} = |P| = 4$   =>   $P$ occurs on pos $10 - 4 - 1 = 5$ of $S$

**runtime:**   Line 1   constant

     2   $n + m$    $\in O(n+m)$

     3   $n + m + 1 \in O(m+n)$

     4   $O(|S|)$    $= O(m+n)$

    5-7   $m + n - 1 - (n+2) = O(m)$

   TOTAL :   $O(m+n)$ time.

**correctness:**

     Since $\$$ not in $P$ and $T$

   $\Rightarrow$   $z_i \leq |P|$   $\forall i$

    if $|z_i| = |P| \Rightarrow$   $S[1 \ldots |P|] = S[i \ldots i + n - 1]$

            "Def       "

            $P$        substring of $T$

        $\Rightarrow$   $P$ in $S$ at pos $i$

         $\equiv P$ in $T$ at pos $n - i - 1$

    if $P$ occurs in $T$ at pos $j$, then

    $P$   occurs in $S$ at pos $i = j + n + 1$

        where $i \in [n+2, m+2]$

         $\Rightarrow |z_i| \geq |P|$

    Since $|z_i| \leq |P|$   $\rightarrow$   $|z_i| = |P|$ &amp;

           occurence at pos $j$

           is reported

             $/\square$

**Thm :**   Simple-exact matching correctly reports

  all occurences of pattern $P$ in text $T$

  in   $O(|P| + |T|)$ time.

There is another important algorithm for pattern matching:

Boyer - Moore Algorithm is standart alg.
that is used in many applications to text search

(e.g. google.. )

BM - Alg. is based on 3 essential ideas:

(1)   Right - to - left scan
(2)   Bad - Character Rule
(3)   Good - Suffix Rule.



[ due to limited time & since we want to cover
further topics, we skip this alg.
more information about this alg in
any Bioinf - textbook. ]

So far:   Z - alg used to preprocess P
& then find P.

Now: preprocess T instead!
( text of static & remains unchanged
( eg DNA ))

# 8. Suffix Trees

classical "real world" problem:

    For given Text & Pattern P
    where and how often does P in Text occur?


Application:
- word processing
- internet search
- Bioinformatics ( T = genom, P = gene sequ.)
- fgrep in Unix
- search for plagiarism
- ⋮

**Def:**
- $\Sigma$ = finite alphabet,
- string $S$ (over $\Sigma$) = sequence of characters in $\Sigma$

Let $S = x_1 x_2 \dots x_\ell$, $x_i \in \Sigma$, $1 \le i \le \ell$

$\qquad |S| = \ell$ = length of $S$

$S[i..j]$ = substring $x_i x_{i+1} \dots x_j$ of $S$, starting at pos $i$, ending at pos $j$

[if $i > j$, then $S[i..j] = \varepsilon$ empty string]

$S[1..j]$ = prefix of $S$
$S[i..\ell]$ = suffix of $S$
prefix/suffix **proper** if it is not $S$ or $\varepsilon$.

$S(i)$ = $i$-th character of $S$

$x, y \in \Sigma$ __match__ if $x = y$, else __mismatch__

$\Sigma^* $ = set of all strings over $\Sigma$, $\Sigma^\ell$ = set of all strings over $\Sigma$ of length $\ell$

$S'$ substring of $S$, if $\exists\ i, j$ st $S' = S[i..j]$

**Aim:** given string $P$ check where/how often $P$ occurs in string $T$

tree data structures can be seen as a collection of
entities (= vertices) that are linked to simulate
a hierarchy
$$= \boxed{\text{rooted tree}}$$
$$= \text{trees } T \text{ where one vertex } \varrho \in V(T) \text{ is called } \underline{\text{root}}$$

Given $\varrho \in V(T)$ we obtain a $\boxed{\textbf{\textit{partial order}}}$ $\leq_T$ on $V(T)$ st
$$x \leq_T y \quad \text{if} \quad y \text{ lies on simple path from } \varrho \text{ to } x.$$
in this case: $x$ is $\boxed{\text{descendant}}$ of $y$ & $y$ is $\boxed{\text{ancestor}}$ of $x$

$v \in V(T)$ is $\boxed{\text{common ancestor}}$ of vertices in $W \subseteq V(T)$, if $w \leq_T v \ \forall w \in W$

a $\leq_T$-minimal common ancestor $v$ of vertices in $W$ is called
$$\boxed{\text{least common ancestor}}$$
$$\boxed{(\text{lca}(W))}$$

$\underline{\text{Exmpl}}$:



the arrows indicate $\leq_T$

Note: $x \leq_T x \quad \forall x \in V(T)$

Notation: $x <_T y$ if $x \leq_T y$ & $x \neq y$

$\underline{\text{Exmpl}}$

$\varrho$

1 has children 2 and x
which are siblings

Leaves are $2, a, b, c$

$\boxed{\text{root}}$
$\boxed{\text{children}}$ of $x$: $y \in V(T)$ st $y <_T x$ & $(xy) \in E(T)$
$\quad x$ is $\underline{\text{parent}}$ of its children

$\boxed{\text{Siblings}}$ are vertices with the same parent
$\boxed{\text{leaf}}$: vertices with no children

<u>Exmpl:</u>    S = honolulu

        S[1..3] = hon ,    prefix & substring of S , no suffix
        S[5..8] = lulu ,   suffix & substring of S , no prefix
        S[5..7] = lul ,    substring of S , no prefix, no suffix


<u>NOTE:</u> empty string ε is substring, prefix, suffix of all strings.


<u>Naive Way</u> :



check  occurence  of P in T
by comparing  1 letter of P with i letter of T
              j letter of P with i + j - 1 letter of T
              n letter of P will i + n - 1 letter of T

                    ∀ i = 1 .. n - m + 1

    ⟹ run time $O(m \cdot n)$


Often,  the text   is  fixed & does not change
         (= long string)
              eg. • collected work of Shakespeare
                  • Genom

To find a pattern P ,   we need a datastructure that represents the text
         (= string)
              so  that  we  can find efficiently P

              To this end :  <u>suffix trees</u>

## Def [suffix tree]   Let s be string of length |S|=m

A suffix tree for S is a rooted tree T (with root $s_r$) that satisfies the following properties:

(X1)  T has precisely m leaves that are uniquely labeled with 1,2...m

(X2)  all inner vertices, (= vertices that are not leaves) except possibly the root, have at least 2 children

(X3)  every edge of T is labeled by a non-empty substring of S

(X4)  For all distinct children $v_1, v_2$ of v we have:  label of edge $(v, v_1)$ & $(v, v_2)$ start with different characters

(X5)  if we concatenate the labels on the edges in order from $s_r$ to leaf $i$, we obtain the suffix $S[i..m]$

Exmpl:   S = a , T   $\overset{s_r}{\underset{1}{\overset{\big|a}{\bullet}}}$   [this is the only case where $s_r$ has only one child, since it |S|>1 &

T. $\overset{}{\underset{m}{\nearrow}}^{S[n..c] \neq \varepsilon}$  we dont get $S[m..m]$, since $|S[m..m]| = 1$ & (X3). ]

S = abc
T



S = abbc



S = aac

What is with $S = aa$ ?          $T \hat{=}$  or 

Observation: not for all strings a suffix tree exists!

Lemma 8.1     IF $S = \underbrace{x_1 \cdots x_k}_{\text{prefix}} \cdots \underbrace{x_i \cdots x_m}_{\text{suffix}}$    st $x_1 \cdots x_k = x_i \cdots x_m$

that is   $S$ contains prefix that is also a suffix of $S$.

Then, no suffix tree for $S$ exists

proof:      Wlog,
we assume that such suffix/prefix are chosen in $S$
to be maximal (there is no longer prefix that is also suffix)

$T$:           or          

lowest common ancestor
of $1$ & $i$   is $\S_T$

Since $x_i = x_1$, this does not
work due to (X4)

lowest common ancestor
of $1$ & $i$   is $v <_T \S_T$

$x_1 \, x_2 \, x_3 \cdots \cdots \, x_k \cdots \cdots x_m$
$\| \quad \| \quad \| \qquad \|$
$x_i \, x_{i+1} \, x_{i+2} \cdots \cdots x_m$

     , $\ell < m$

$\lightning$ (x4) since
$x_{i+\ell} = x_{\ell+1}$

$\Rightarrow$ only possibility:      $\lightning$ (x3)

$\Rightarrow$ if suffix of $S$ is also prefix of $S$ no suffix tree exists!

IDEA: simply add special symbol $\$$ at end of $S$,
where $\$$ does not occur in $S$.

From here on: $S = x_1 \cdots x_m$ with $x_m = \$$

Example: $S = xab\,xa\,\$$ $\qquad$ ( for "$xabxa$" no suffix tree)
$\phantom{Example: S = }1\ 2\ 3\ \ 4\ 5\ 6$



Def.

Let $v \leq_T w$ & $P$ be the unique simple path from $v$ to $w$



label of $P$ = concatenation of labels of edges in order from $v$ to $w$.

By slight abuse of notation, we also consider simple paths in suffix tree that may end on edge (and not a vertex)



$\alpha$ = label of path from $v$ to $w$

$(w\,w') \in E(T)$
with label $t_1 \cdots t_e$

$P$ ends here

label of $P = \alpha\, t_1 \cdots t_i$

# How to construct suffixtree & how does this help?

## Alg. SuffixTree (Idea)

Input:  $S = x_1 \ldots x_m$

1) Construct $T_1 = $



2) Assume $T_i$ is constructed, then construct $T_{i+1}$
   $(i < m)$
   as follows:

(a) Find path $P$ in $T_i$ that starts in $\S_{T_i}$ with __longest label__ that is a __prefix__ of
   $S[i+1 \ldots m]$    [path could end in $\S$]

    (i) By similar arguments as in proof of L. 8.1
   __this path is uniquely determined__
   since no two edges $(v, v_1), (v, v_2)$, $v_1, v_2$ children of $v$
   have labels starting with same symbol $(X4)$.

    (ii) __this path will never end in leaf of $T_i$__
   since $P$ starts at root & concatenation of edge labels
   to leaf $i$    yields $S[i \ldots m]$ where $|S[i \ldots m]| > |S[i+1 \ldots m]|$
   [by induction - Exnc]

(b) This path $P$ either ends in edge or non-leaf vertex of $T_i$

$P$ ends in vertex:



w no leaf or w has children.

$\alpha$ prefix of $S[i+1 \ldots m]$, $m = \$$
But $\alpha \neq S[i+1 \ldots m]$
otherwise $S[i+1 \ldots m] = S[\ell \ldots m]$
    for some $\ell$ but $\ell < i+1$ ⚡
$\Rightarrow \alpha = S[i+1 \ldots m-j]$, $j > 1$

$\Rightarrow$ add __edge $(w, i+1)$__
   with label $\beta = S[m-j+1 \ldots m]$

P ends in edge :



$\alpha = S[i+1 \ldots m-j]$

label edge $(w,w') = t_1 \ldots t_r \, t_{r+1} \ldots t_s$

$\Rightarrow \quad \alpha = \alpha' t_1 \ldots t_r \quad$ for some $1 \le r < s$

add **new vertex v**
on edge $(w, w')$
& put $label(wv) = t_1 \ldots t_r$
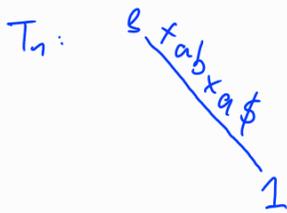$\quad\quad label(vw') = t_{r+1} \ldots t_s$



+ add **new edge $(v, i+1)$**
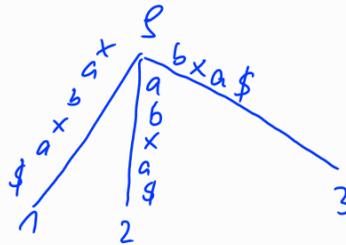with $label(v, i+1) = S[m-j+1 \ldots m] = \beta$

For both cases, $T_{i+1}$ contains now path from $s_{T_{i+1}}$ to $i+1$
with label $S[i+1 \ldots m]$

# Exmpl: $\quad S = xabxa\$$

$T_1:$



$T_2 / T_3 \quad$ (longest path ends in $s$ )



$T_4: \quad S[4..m] = xa\$ \qquad$ longest path P



$T_5: \quad S[5..m] = a\$$

**SUFFIX-TREE(S)**

init $T$ as  and label $(g_T, 1) = S[1..m]$

FOR ( $i = 2..m$ ) DO

  $(end, i', \ell') = $ FIND-LONGEST-PATH $(T, S[i..m])$   // returns end of path, that is, either $end = v$ or $end = e$ & index $i'$ & $\ell'$

  IF ( $end = $ edge $e$ )   // $e = (uw)$ with label $g$

    remove $e = (uw)$
    add new vertex $v$
    add new edges $e_1 = (uv)$, $e_2 = (v, w)$
      & label $e_1 = g[1..\ell']$
      label $e_2 = g[\ell'+1 .. |g|]$

  add edge $(v, i)$ to $T$ with label $S[i+i' .. m]$

---

**FIND-LONGEST-PATH $(T, S')$**    <span style="color:red">[Sketch → more details in Appendix]</span>

"Follow path from $g_T$ to $v/e$ as long as possible, i.e, as long as letter on this path match with letters $S' = x_{\bar{i}}.. x_m$"

& return corresponding positions $i', \ell'$
    + if end is edge or vertex.

---

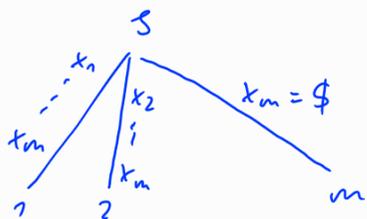**Theorem [UKKONEN]:**    for $S$ of length $m$ the suffix tree can

[without proof]    be constructed in $O(m)$ time

<span style="color:red">[quite sophisticated pointer-adjustments].</span>

# Space complexity

Suffix tree without labels $O(m)$
but we need to store labels!

worst case,



each edge label of $(\beta, i)$ is of size $i$

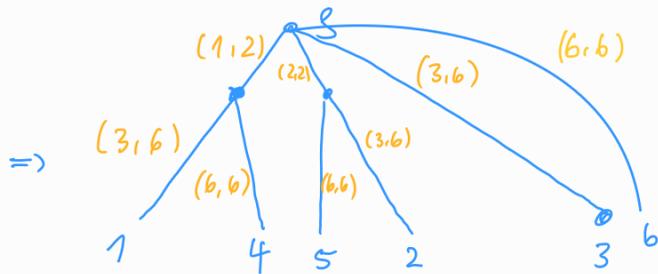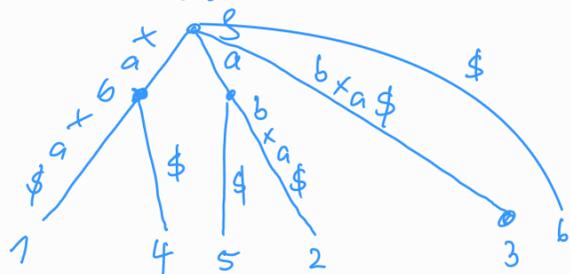$\Rightarrow \sum_{i=1}^{m} i = O(m^2)$ space

$O(m^2)$ space is bad   (eg. genome of small bacteria has
$10^6$ characters $\sim$ 1 TB storage

IDEA:

Instead of saving label $S[i..j]$ for $c$
we only save pair $(i,j)$

$=$ compressed suffix tree

$\Rightarrow O(m)$ space ( same space as text $S$)

Exmpl: $S = x\,a\,b\,x\,a\,\$$
          $1\ 2\ 3\ 4\ 5\ 6$



# Why Suffix trees?

A: Examples!

In what follows, we assume to
have Ukkonen's version: $O(m)$ time

**"etact text-search"**

Given (long) text S & palter P (=string)

Does P in S occur?

Let T be suffix tree of S.


Call   FLP (T, P)   once   O(|P|) time

return value is   (end, index i, ...)

st   P[1..i'] corresponds to label of path in T

that starts in $\int T$


Observe, any path in T von $\int$ to leaf i

corresponds to   S[i..m]

&   P[1..i'] corresponds therefore to a
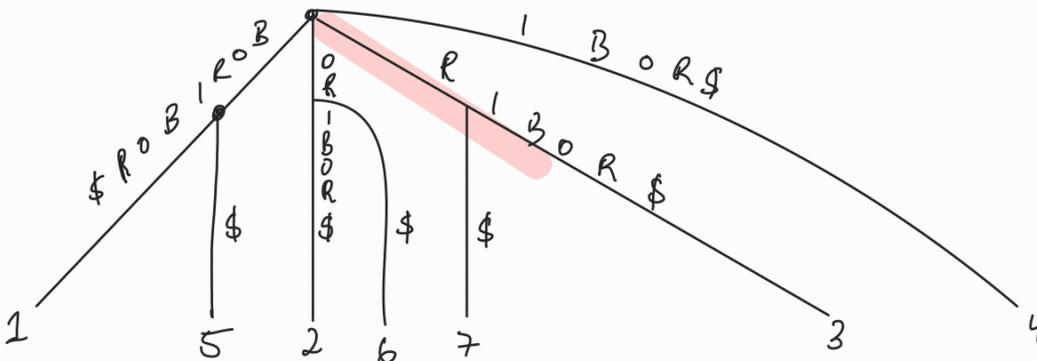
prefix of some suffix of S, that is,

a substring of S.


=)   if  i' = |P|   => P occurs in T

else, not.

=)   runtime  O(|P|)

Exmpl:   T= BORIBOR$.     P = RIB



[every exit letter appears once at first letter of edge ($V$) for some u.]

Exmpl 2   "substring - database - search"

Given strings  $S_1 \ldots S_\ell$   (Database of Texts)
Does  P  occur  in (at least) one of the $S_i$?

=> let  $\$_1 \ldots \$_\ell$  pairw. distinct symbols that do not occur in
                                                          any $S_i$ & P

put  $S = S_1 \$_1 S_2 \$_2 \cdots S_\ell \$_\ell$
Construct suffix tree  T  for  S    ( $O(|S_1| + |S_2| + \ldots |S_\ell|)$ time )
use IDEA of Exmpl 1


Exmpl 3   Find all occurences of P

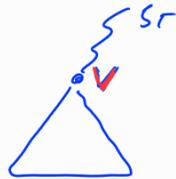given  S  with suffix tree  T  & pattern P with $|P| = n$

Find all occurences of  P  in  T, that is,
all indices  i  st  $P = S[i \ldots i+n-1]$


call  FLP( T, P )
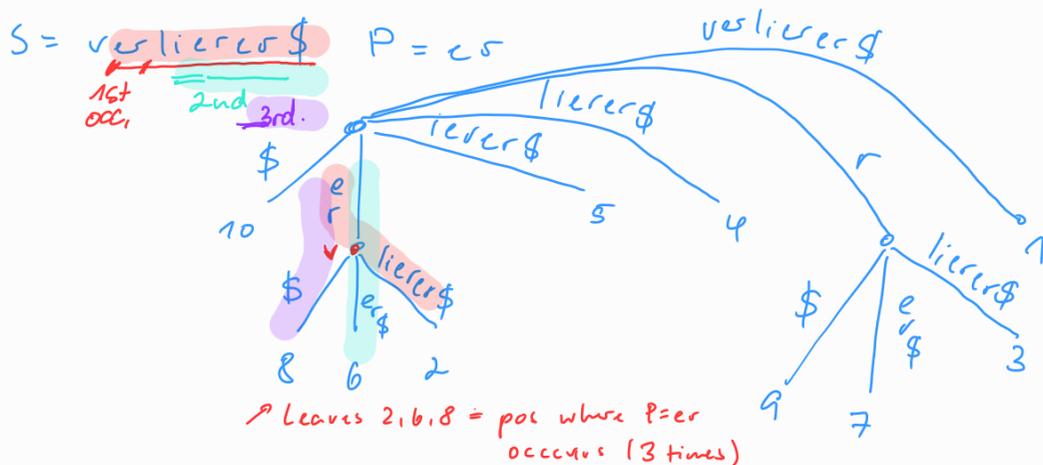      $\longrightarrow$  this gives  end = v  or  end = edge ($u$✔)



# leaves $i \leq_T v$ = # occurrences of P in S

& P occurs in pos i $\forall$ all leaves $i \leq_T v$.

[ Exercise :  works in  $O(|P| + k)$ time )
                                              ↗
                                         k = nr of
                                          occurences
                                           of P

$S = $ verlierer$\$$   $P = er$        verlierer$\$$



Leaves 2,6,8 = pos where P=er
                occurs (3 times)

_Exmpl 4_        ==Find longest substring in S that occurs at least 2 times==

[Application: Find repeated regions in genom]

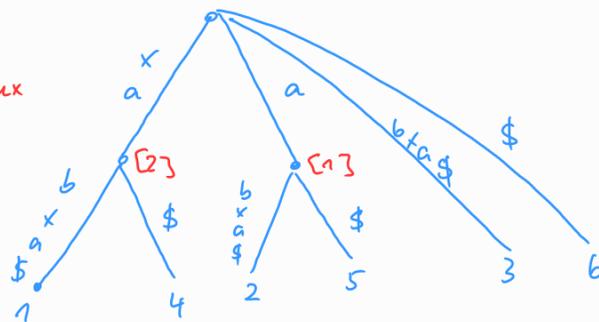M I S S I S S I P P I        i S S i

string depth of $v$ in $T$
is $|\alpha|$

$S_T$
$\alpha$ = path label of path from $S_T$
to $v$
$v$

$S =$ verliercr $\$$

string depth
of inner vertex

verliercr$\$$

liercr$\$$

lercr$\$$

r

$\$$        $e$
$r$
10        [2]
liercr$\$$
$\$$  $er$\$$
8   6   2

$S$   4   [1]
liercr$\$$
$\$$  $e$
$\$$
9  7   3

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

string depth
of inner vertex

$x$
$a$            $a$
$b$ $a$ $\$$     $\$$
[2]           [1]
$b$    $\$$    $b$    $\$$
$x$          $x$
$a$          $a$
$\$$          $\$$
$\$$           5     3   6
1     4   2

● every inner vertex has 2 children &   first symbol on edges with
        (at least)                      same parent have distinct symbols

⟹ length of longest substring that occurs at least 2 times

= maximum string-depth of inner vertices.

[Exus. : works in $O(|S|$ time.]

efficient detection of pairwise overlaps

given $\mathcal{J} = \{S_1 \cdots S_N\}$ set of strings ( e.g. sequenced DNA fragments)

Aim: define $ov(S_i, S_j)\ \forall\ i,j$
$\qquad\qquad$ " size of largest overlap

$S = S_1\, \$_1\, S_2\, \$_2 \cdots S_N\, \$_N$

$S_i\ \overline{\phantom{xxxx}}$

$\underline{\phantom{xxx}}\ S_j$

$S_1 = ABA \qquad\longrightarrow\qquad ov(S_1 S_2) = 2$

$S_2 = BAA \qquad\qquad\qquad ov(S_2 S_1) = 1$

$$S = \overset{1}{A}\ \overset{2}{B}\overset{3}{A}\ \overset{4}{\$_1}\ \ \overset{5}{B}\overset{6}{A}\overset{7}{A}\ \overset{8}{\$_2}$$

name of vertex



$P_1$ = path starting with label $S_1 \$_1$

$P_2$ = path starting with label $S_2 \$_2$

for non-leaf $x$: $\quad L(x) = \{\, i \mid \exists\ \text{leaf } v \text{ s.t. } (xv) \text{ edge in suffix tree}$
$\qquad\qquad\qquad\qquad \& \text{ label } (xv) \text{ starts with } \$_i \,\}$.

string dept of $x$ = # characters on path $\$$ to $x$.

For (j=1 .. N)

    x = ƒ (root)

    WHILE (x not leaf)

        FOR (all $i \in L(x)$, $i \neq j$)

            IF (depth(x) < min ($|S_i|$, $|S_j|$))

                or $(S_i, S_j)$ = depth (x)

        x ← child of x   on $\tau_{ij}$

---

$j=1$, $x = ƒ$

FOR ($i \in L(ƒ) = \{1,2\}$, $i \neq 1$)

$\Rightarrow$ $i = 2$

   depth(v) = 0 < min($|S_1|$, $|S_2|$)

   $\Rightarrow$  or $(S_2, S_1) = 0$

x ← u

FOR ($i \in L(u) = \{1\}$, $i \neq 1$)

        nothing

x ← leaf  stop.

$\bar{j}=2$, $x - ƒ$

For ($i \in L(ƒ) = \{1,2\}$, $\bar{i} \neq 2$)

$\Rightarrow$ $i = 1$

   dept (ƒ) = 0 < min ($|S_1|$, $|S_2|$)

   or $(S_1, S_2) = 0$

x ← w

For ($i \in L(w) = \{1\}$, $\bar{i} \neq 2$)

   $\Rightarrow$ $i = 1$

   depth (w) = 2 < min ($S_1$, $S_2$)

        $\Rightarrow$ or $(S_1, S_2) = 2$

x ← leaf  stop.

---

[ without proof of correctness   or runtime analysis ]

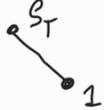$$O(N \, \|S\|)$$

Exercise  / Book "Alg. Aspects of Bioinf.")

# Appendix

## Details of $O(m^2)$ Algorithm

## SUFFIX-TREE (S)

init $T$ as



and label $(\S_T, 1) = S[1..m]$

FOR( $i = 2 .. m$ ) DO

$\quad$ $(end, i', l') = $ FIND-LONGEST-PATH $(T, S[i..m])$ $\quad$ // returns end of path, that is, either $end = v$ or $end = e$ & index $i'$ & $l'$

$\quad$ IF ( $end = $ edge $e$ ) $\quad$ // $e = (uw)$ with label $\gamma$

$\qquad$ remove $e = (uw)$
$\qquad$ add new vertex $v$
$\qquad$ add new edges $e_1 = (uv)$, $e_2 = (v, w)$
$\qquad\qquad$ & label $e_1 = \gamma[1..l']$
$\qquad\qquad$ label $e_2 = \gamma[l'+1..|\gamma|]$



$\quad$ add edge $(v, i)$ to $T$ with label $S[i + i' .. m]$

## FIND-LONGEST-PATH $(T, S')$

$\quad j = 1, \quad v = \S_T$

$\quad$ WHILE ( $j \leq |S'|$ ) DO

$\qquad$ Find edge $e = (vw)$ in $T$, $w \leq_T v$, whose label starts with $S'(j)$

$\qquad$ IF ( such edge does not exist )

$\qquad\qquad$ RETURN $(v, j-1, \emptyset)$ // path ends in $v$

$\qquad$ Let $\gamma$ be label of $e$

$\qquad$ $l = 1$

$\qquad$ WHILE ( $j \leq |S'|$ & $l \leq |\gamma|$ & $S'(j) = \gamma(l)$ ) DO

$\qquad\qquad$ $j = j + 1$
$\qquad\qquad$ $l = l + 1$

$\qquad$ IF ( $l \leq |\gamma|$ ) $\quad$ // while loop ended at some point before "$S'(j) = \gamma(l)$ $\forall l = 1 .. |\gamma|$"

$\qquad\qquad$ RETURN $(e, j-1, l-1)$ $\qquad$ ie, $S(j) \neq S(e)$

$\qquad$ $v = w$ $\quad$ // go to consider edges starting at $w$
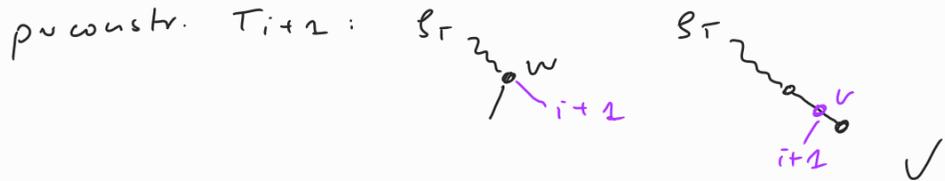
$\quad$ RETURN $(v, |S'|, \emptyset)$

<u>Prop. 8.2</u>   Alg. SUFFIXTREE correctly computes suffix tree
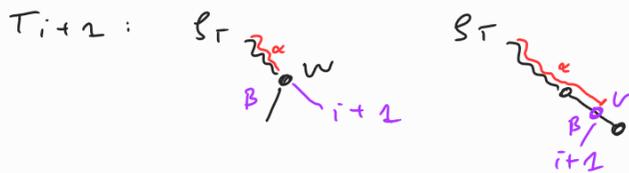            for $S = x_1 \cdots x_m$ , $x_m = \$$ .


<u>proof:</u>   (X1):   $T_1$ has 1 leaf, $T_2$ has 2 leaves ... $\overset{ind.}{\ldots}$ $T_m$ has $m$ leaves ✓
(sketch)   (X2):   $T_1$ :   ✓

            Assume, for $T_i$, each inner vertex (except possibly $\$$)
                        has at least 2 children

            pre constr. $T_{i+1}$ :    ✓


    (X3)   by construction & induction   each edge
            has non-empty string as label


(X4)   $T_1$ ✓   Assume for $T_i$ $\forall v$ : label $(vv_1), (vv_2)$ starts with different
                                                symbol, where $v_1, v_2$ are
                                                        children of $v$.

        $T_{i+1}$ :   

        Since $\alpha$ is <u>longest</u> substring that is prefix of

        $S[i+1 \ldots m] = \underbrace{x_{i+1} \cdots}_{\alpha} \underbrace{x_r \cdots x_m}_{\beta}$        $\beta(1) \neq x_{r+1}$ as otherwise
                                                                    $\alpha$ not longest prefix !
                                                                        ✓

(X5)   by constr. & induct.   concatenating edge-labels from $\$$ to $i$
                                            yields $S[i \ldots m]$.

                                                            $\diagup \square$

[suffix-
tree has $O(m)$ edges & vertices]

Runtime: SUFFIXTREE(S) , $S = x_n \ldots x_m$

add/remove etc in $O(1)$ time

$\rightarrow$ m times FIND-LONGEST-PATH (FLP) is called.

FPL: • in each of the $|S'|$-steps in 1st while-loop

Find-edge $\rightarrow$ for $v$ all neighbors $w$ are considered

$\Rightarrow$ $\deg_T(v)$ many vertices

Over all calls of 1st-while loop, "FIND-edge" is called,

$$\leq \sum_{v \in V} \deg(v) = 2|E| = O(|E|) \text{ times}$$
$$= O(m)$$

• 1st + 2nd while loop:

WHILE $(j \leq |S'|)$ DO
: 
  WHILE $(j \leq |S'| \,\&\, \ell \leq |x| \,\&\, S'(j) = x(\ell))$ DO
  | $j = j+1$
  _ :

$\Rightarrow$ $j \leq m$ $\Rightarrow$ $O(m)$ time

• Remaining operation in 1st-while loop: $O(1)$

$\Rightarrow$ FLP runs in $O(m)$ time

$\Rightarrow$ SUFFIX-TREE(S) runs in $O(m^2)$ time.

---

Theorem [Ukkonen]: for $S$ of length $m$ the suffix tree can
[without proof] be constructed in $O(m)$ time