

3. Exact Paterson Matching

String Matching

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

Aim: find occurrences of P in T

Exmpl: $T =$ ^{1 2 3 4 5 6 7 8 9 10 ...} AATGCA**T**GCA
 $P =$ AT**G**

P occurs on position 2, 6, ...

This has applications in:

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

Basics:

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

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

$$S(i) = x_i$$

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

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

S: ^{1 2 3 4 5 6 7 8} HONOLULU
HONOLU LU
 $S[1..4]$ prefix
 $S[7..8]$ suffix
 $S(S) = L$

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

Naive Method

NAIVE (P, T)

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

Occurrences = \emptyset

FOR ($i = 1 \dots |T| - |P| + 1$)

// loop over "pos." in T

FOR ($j = 1 \dots |P|$)

// 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 = 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: $|P| \cdot (|T| - |P| + 1)$ (worst case)

See example above: $2 \cdot 5 = 10$ comparisons

Q: What is least nr of character comparisons?

A: $|T| - |P| + 1$ (best case)

T = A A A A A A

P = B A

⇒ RUNTIME NAIVE-alg

$|T| - |P| + 1$
 $|P|$

1. loop

2 loop

$$\Rightarrow |T|(|P| - |P|^2 + 1) \leq |T||P|$$

$|T| \geq |P|$

⇒ $O(|T||P|)$ 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)
& $i > 1$.

$z_i := z_i(S) =$ length of longest substring
in S that
starts at position i &
matches prefix of S

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
z_i	-	1	0	0	3	1	0	0	2	1	0

DEF (z-Box): $\forall i > 1$ s.t. $z_i > 0$, the
z-Box at i is
the interval $[i, i + z_i - 1]$

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
z_i	-	1	0	0	3	1	0	0	2	1	0

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

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

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

(Red brackets in the original image group the 'a's at positions 2, 5-6, and 9-10.)

These intervals $[i, j]$ correspond
to $S[i..j]$ = longest substring
starting at $i \geq 1$ &
matches prefix of S

DEF

$\forall i > 1 :$

(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 \uparrow

(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

Formal: $r_i = \max_{2 \leq j \leq i} \{ j + z_j - 1 : z_j > 0 \}$

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

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 EXPLICIT character comparisons
as much as possible.

Overview of steps of z -alg

① compute z_2 (explicit comparison of $S[1..n]$ with $S[2..n]$ until mismatch)

IF ($z_2 > 0$)
ELSE

put $l_2 = 2, r_2 = z_2 + 1$
put $l_2 = r_2 = 0$

②

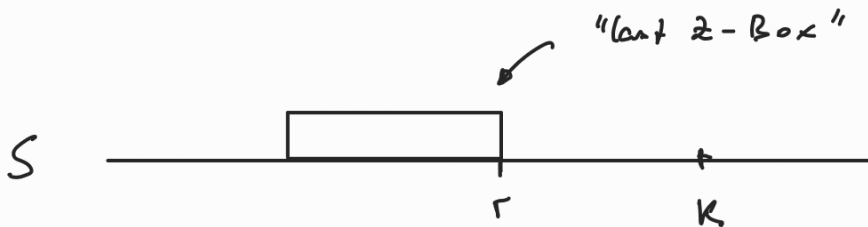
k -th step ($k \geq 2$).

\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: $r < k$



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

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

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			

$k=9 \rightarrow r_g = 8 < k=9$ need to compare $S[9..11]$ with $S[1..n]$

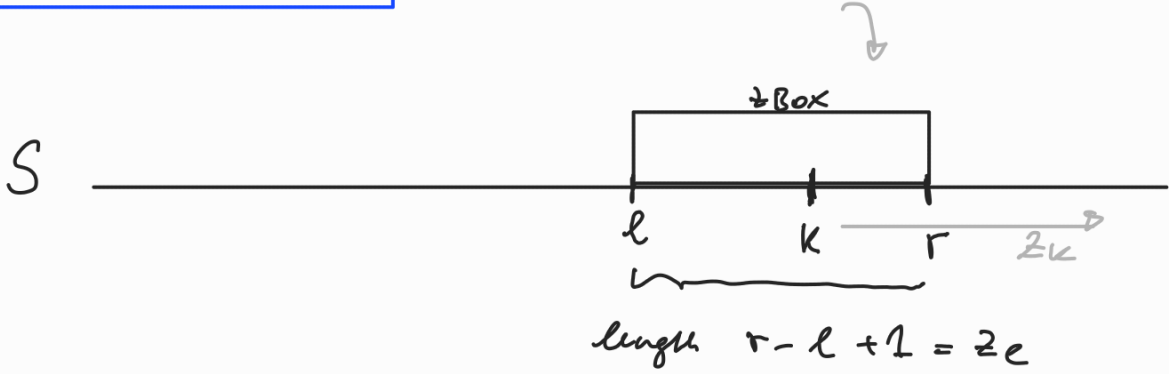
$\Rightarrow S[9] = S[1] \checkmark$
 $S[10] = S[2] \checkmark$
 $S[11] \neq S[3] \times$

$\Rightarrow z_g = 2, l_g = 9$

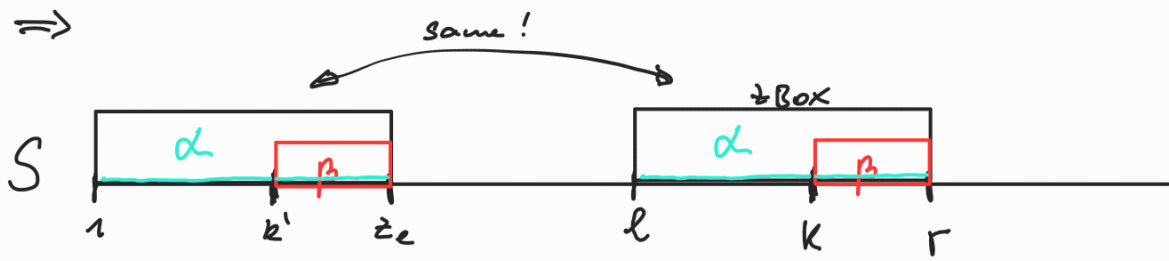
$r_g = 9 + 1 - 1 = 10$

Case 2: $r \geq k$

[Now we have some pre-knowledge about about z_k]



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



$$S[1..z_c] = S[l..r] = \alpha$$

$$S[k'..z_c] = S[k..r] = \beta$$

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

Question: what we know about β ?

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

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

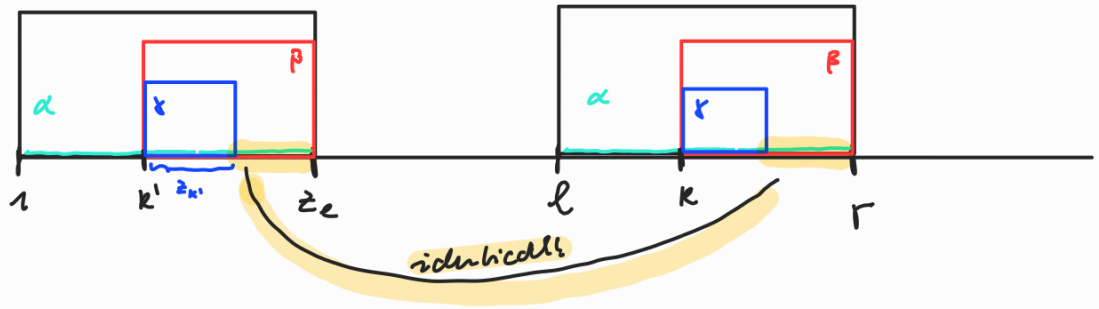
$k' = k - l + 1$

and leads to following subcases.

2A, 2B

2A

$$z_{k'} < |\beta| = r - k + 1$$



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

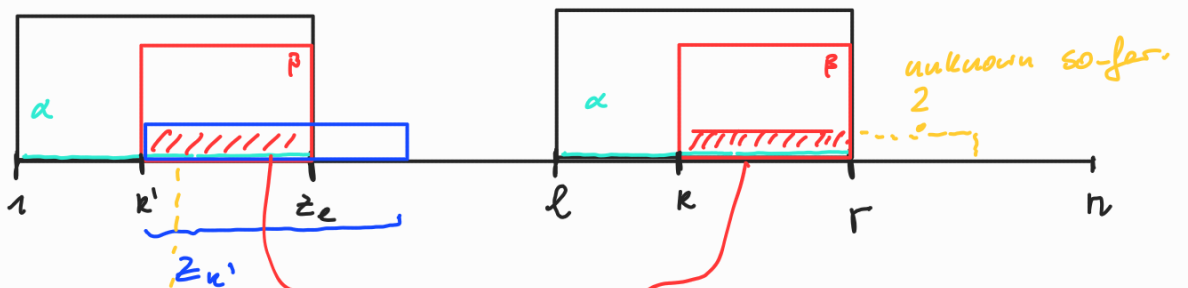
Since $\beta = \beta$
 left right $\Rightarrow \gamma$ (right) is longest substring starting at k & is prefix of S .

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

l & r remain unchanged.

2B

$$z_{k'} \geq |\beta| = r - k + 1$$



$\beta + 1$
 pos.

this part matches already,
 i.e. β (right) is part of longest
 substring starting at k & prefix
 of S .

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

$$\text{so } z_k \geq |\beta| = r - k + 1$$

However $Z_k > |S|$ may be possible.

\Rightarrow compare $S[r+1..n]$ ^[yellow part above]
 with $S[\beta|+1, \dots, n]$ until mismatch
 //
 $r - k + 1 + 1$
 $- r - k + 2$
 As possibly extended
 Z-Box starting at k

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

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

Z Algorithm

```

1:  $r \leftarrow l \leftarrow 0$ 
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  $r \leftarrow k + Z_k - 1, l \leftarrow k$ .
9: // Case 2:
10:  else
11:     $k' \leftarrow k - l + 1$ 
12: // Case 2a:
13:   if  $Z_{k'} < r - k + 1$  then
14:      $Z_k \leftarrow Z_{k'}$ 
15: // Case 2b:
16:   else
17:      $l \leftarrow k$ 
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:    $Z_k \leftarrow q - k, r \leftarrow q - 1$ 

```

later discussion implies:

Thm: Z -alg. correctly computes all
 $Z + \text{Box } Z_i, i > 2$

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 2-Alg we either use 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].

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

$$\Leftrightarrow r_k > r_{k-1} + m_k$$

$$\Leftrightarrow r_k - r_{k-1} > m_k \quad \checkmark$$

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

(implicit in algo) $\Leftrightarrow r_k - r_{k-1} = m_k = 0. \quad \checkmark$

Now Case 2A/2B.

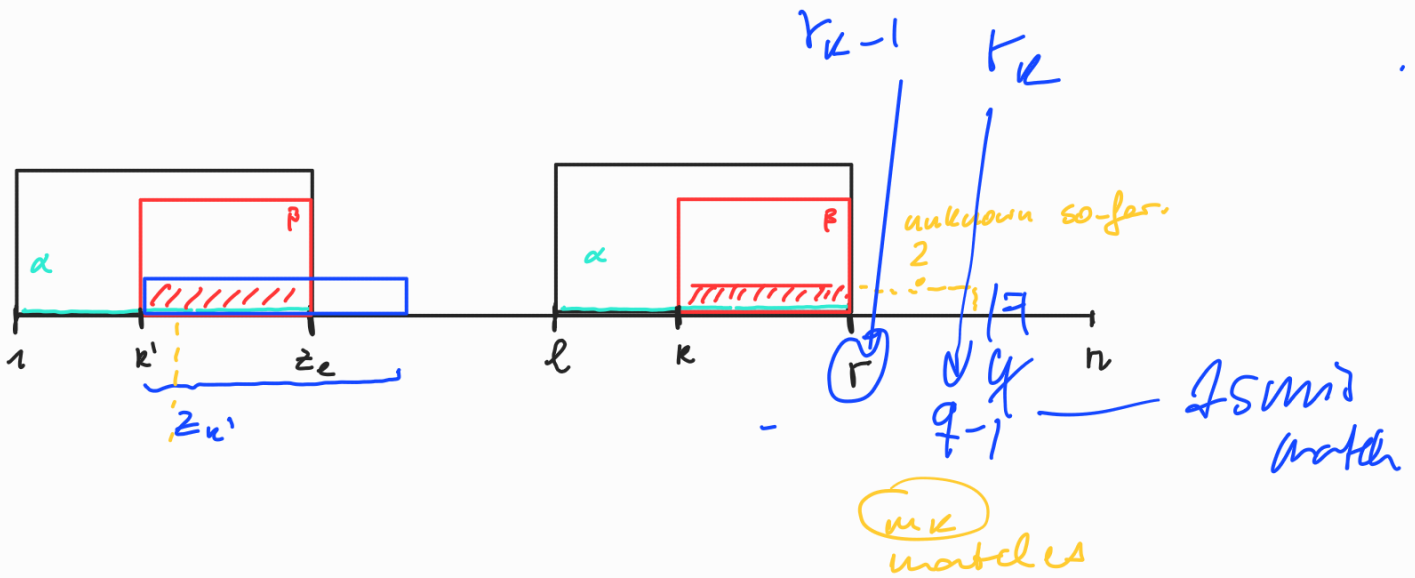
CASE 2A: in Alg $r = r_{k-1}$ remains unchanged,

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

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

$\Leftrightarrow r_k - r_{k-1} = m_k \checkmark$

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



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

$$n := |S| \Rightarrow \sum_{k=2}^n c_k \leq \sum_{k=2}^n 1 + \sum_{k=2}^n m_k = n-1 + \sum_{k=2}^n m_k$$

(miss.) (match)

$\leq n-1 + (r_2-r_1) + (r_3-r_2) + (r_4-r_3) \dots + (r_n-r_{n-1}) = n-1 + \underbrace{r_n-r_1}_{\leq n} \leq 2n-1 \in O(n) \quad \square$

(claim)

\downarrow
in Alg $r_1 = l_1 = 0$

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	-									

$$\text{init } r = l = 0$$

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

compare $S[2 \dots n]$ with $S[1 \dots n]$

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

$\Rightarrow z_k = 0$, 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								

$$k=3, \quad k=3 > r=0 \quad (\text{Case 1})$$

compare $S[3 \dots n]$ with $S[1 \dots n]$

$$S(3) = S(1)$$

$$S(4) = S(2)$$

$$S(5) = S(3)$$

$$S(6) \neq S(4)$$

\Rightarrow

$$z_k = 3 > 0$$

$$r = 3 + 3 - 1 = 5$$

$$l = 3$$

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							
l_i	-	0	3							
r_i	-	0	5							

$$k=4, \quad r=5 \geq k=4 \quad (\text{Case 2})$$

$$(\alpha = aca, \beta = ca)$$

$$k^1 = k - l + 1 = 4 - 3 + 1 = 2$$

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

$$\Rightarrow \text{Case 2A: } z_4 = z_{k^1} = 0, \quad r, l \text{ unchanged.}$$

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						
l_i	-	0	3	3						
r_i	-	0	5	5						

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

$$k^1 = k - l + 1 = 5 - 3 + 1 = 3$$

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

$$\Rightarrow \text{Case 2A again and so on...}$$

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 $T \cup P$

- 1 occurrences = \emptyset
- 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	<u>bb</u>	<u>aca</u>			
z_i	(doesn't matter...)					0	3	1	0	<u>4</u>	1	0	0	0

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

Runtime:

Line 1	constant
2	$n + m \in O(n + m)$
3	$n + m + 1 \in O(n + m)$
4	$O(S) = O(n + m)$
5-7	$m + n - 1 - (n + 2) = O(m)$

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

correctness:

Since $\$$ not in P and T
 $\Rightarrow z_i \in |P| \forall i$

if $|z_i| = |P| \Rightarrow$ $S[1 \dots |P|] \stackrel{\text{"pat } P}{=} S[i \dots i + n - 1] \stackrel{\text{"substring of } T}{=}$
 $\Rightarrow P$ in S at pos i
 $= P$ 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 $|z_i| \leq |P| \rightarrow |z_i| = |P|$ &
 occurrence at pos j
 is reported

□

Then: Simple-exact matching correctly reports
 all occurrences of pattern P in text T
 in $O(|P| + |T|)$ time.

[Back to slides]

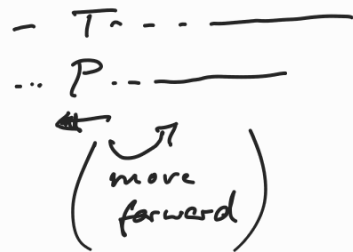
There is another important algorithm for pattern matching:

Boyer-Moore Algorithm is standard alg.

that is used in many applications for 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 Brief-textbook.]

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

Now: preprocess T instead!
(text of static & remains unchanged (e.g. 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 ($T = \text{genom}$, $P = \text{gene sequ.}$)
 - `fgrep` in Unix
 - search for plagiarism
 -

Def:

- Σ = finite alphabet,
- string S (over Σ) = sequence of characters in Σ

Let $S = x_1 x_2 \dots x_l$, $x_i \in \Sigma$, $1 \leq i \leq l$

$|S| = l$ = 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 $S[i..j] = \epsilon$ empty string]

$S[1..j]$ = prefix of S

$S[i..l]$ = suffix of S

prefix/suffix proper if it is not S or ϵ .

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

$x, y \in \Sigma$ match if $x = y$, else mismatch

Σ^* = set of all strings over Σ , Σ^l = set of all strings over Σ of length l

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
 = rooted tree

= trees T where one vertex $f \in V(T)$ is called root

Given $f \in V(T)$ we obtain a partial order \leq_T on $V(T)$ st

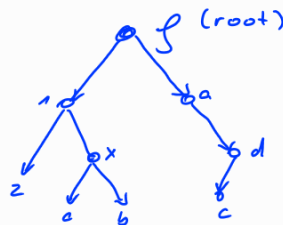
$x \leq_T y$ if y lies on simple path from f to x .

in this case: x is descendant of y & y is ancestor of x

$v \in V(T)$ is 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 least common ancestor ($\text{lca}(W)$)

Exmpl:



the arrows indicate \leq_T

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

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

Exmpl

f

1 has children 2 and x which are siblings

leaves are 2, a, b, c

root

children of x : $y \in V(T)$ st $y <_T x$ & $(xy) \in E(T)$

x is parent of its children

Siblings are vertices with the same parent

leaf: vertices with no children

Exmpl: $S = \text{honolulu}$

$S[1..3] = \text{hon}$, prefix & substring of S , no suffix

$S[5..8] = \text{lulu}$, suffix & substring of S , no prefix

$S[5..7] = \text{lul}$, substring of S , no prefix, no suffix

NOTE: empty string ϵ is substring, prefix, suffix of all strings.

Naive way:



check occurrence of P in T

by comparing i letter of P with i letter of T
 j letter of P with $i+j-1$ letter of T
 n letter of P with $i+n-1$ letter of T

$\forall i = 1 \dots n-m+1$

\Rightarrow run time $O(n^2)$

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: suffix trees

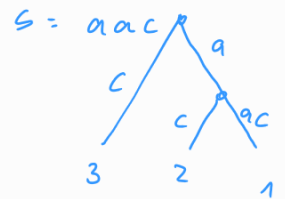
Def [suffix tree] let S be string of length $|S|=m$

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

- (X1) T has precisely m leaves that are uniquely labeled with $1, 2, \dots, m$
- (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 r to leaf i , we obtain the suffix $S[i..m]$

Exmpl: $S = a$, T 

[this is the only case where r has only one child, since if $|S| > 1$ & $T: [S[1..1]] \neq \epsilon$ we don't get $S[m..m]$, since $[S[m..m]] (= 1)$ & (X3).]



How to construct suffix tree & how does this help?

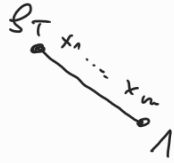
Alg. SUFFIX TREE (Idea)

// iteratively build suffix trees

$T_1 \dots T_m$ by adding $S[1..m], \dots, S[m..m]$
in previously constructed T_i

Input: $S = x_1 \dots x_m$

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 ST_i
with longest label that is a prefix of
 $S[i+1..m]$ [path could end in $\{\}$]

(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 ($\neq \emptyset$).

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

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

P ends in vertex:

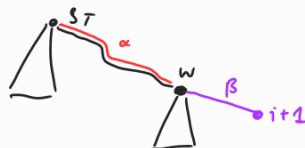


w no leaf $\Rightarrow w$ has children.

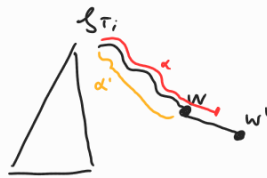
α prefix of $S[i+1..m]$, $m = j$
But $\alpha \neq S[i+1..m]$
otherwise $S[i+1..m] = S[l..m]$
for some l but $l < i+1$

$\Rightarrow \alpha = S[i+1..m-j]$, $j > 1$

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



P ends in edge:



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

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

$$\Rightarrow \alpha = \alpha' t_1 \dots t_r \text{ for some } 1 \leq r < s$$

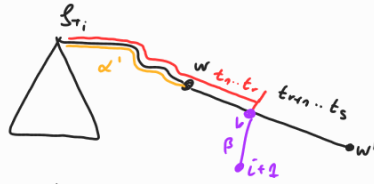
add new vertex v
on edge (w, w')

& put $\text{label}(wv) = t_1 \dots t_r$

$$\text{label}(vw') = t_{r+1} \dots t_s$$

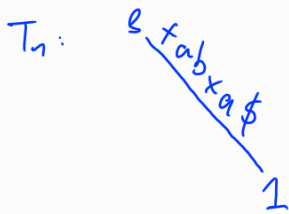
+ add new edge $(v, i+2)$

$$\text{with label}(v, i+2) = S[m-j+2 \dots m] = \beta$$

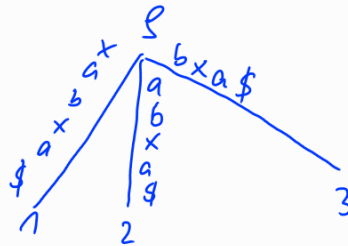


For both cases, T_{i+2} contains now path from s_{i+1} to $i+2$
with label $S[i+2 \dots m]$

Exmpl: $S = xabxa\$$

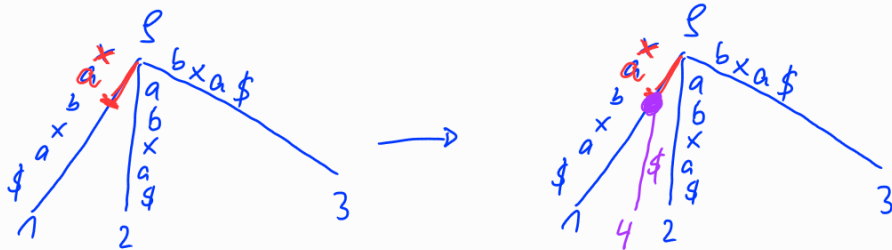


T_2/T_3 (longest path ends in $\$$)

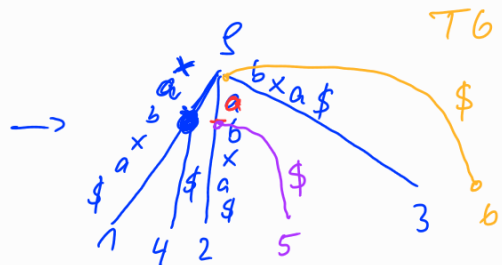
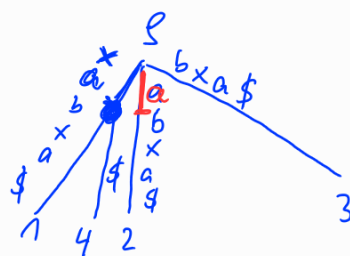


T_4 : $S[4..m] = xa\$$

longest path P



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



SUFFIX-TREE(S)

init \mathcal{T} as  and label $(s_{\mathcal{T}}, 1) = S[1..m]$

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

(end, i' , l') = FIND-LONGEST-PATH (\mathcal{T} , $S[i..m]$)

// returns end of path, that is,
// either end = v or end = e
// & index i' & l'

IF (end = edge e) // $e = (uv)$ with label γ

remove $e = (uv)$
add new vertex v

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



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

FIND-LONGEST-PATH (\mathcal{T} , s')

[Sketch \rightarrow more details in Appendix]

"Follow path from $s_{\mathcal{T}}$ to v/c as long as possible, i.e., as long as letters on this path match with letters $s' = x_{i'}..x_m$ "

& return corresponding positions i' , l'

+ if end is edge or vertex.

Theorem [Ukkonen]:

[without proof]

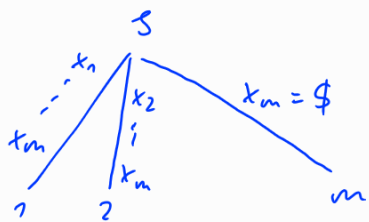
for S of length m the suffix tree can be constructed in $O(m)$ time

[quite sophisticated pointer-adjustments].

Space complexity

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

worst case,



each edge label of (ℓ, i) is of size i
 $\Rightarrow \sum_{i=1}^m i = O(m^2)$ space

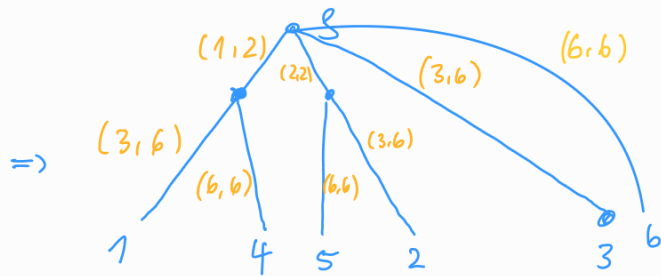
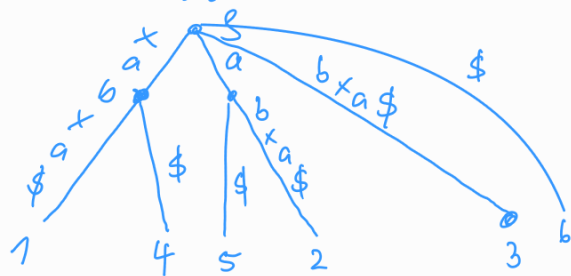
$O(m^2)$ space is bad (eg. genome of small bacteria has 10^6 characters \sim 1TB storage)

IDEA:

instead of saving label $S[i..j]$ for e
 we only save pair (i, j)
 = **compressed suffix tree**

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

Exmpl: $S = xabxa\$$
 1 2 3 4 5 6



Why suffix trees?

A: Examples!

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

Exmpl 1 "exact text-search"

Given (long) text S & pattern 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 $\$T$

Observe, any path in T from $\$$ to leaf i corresponds to $S[i..m]$

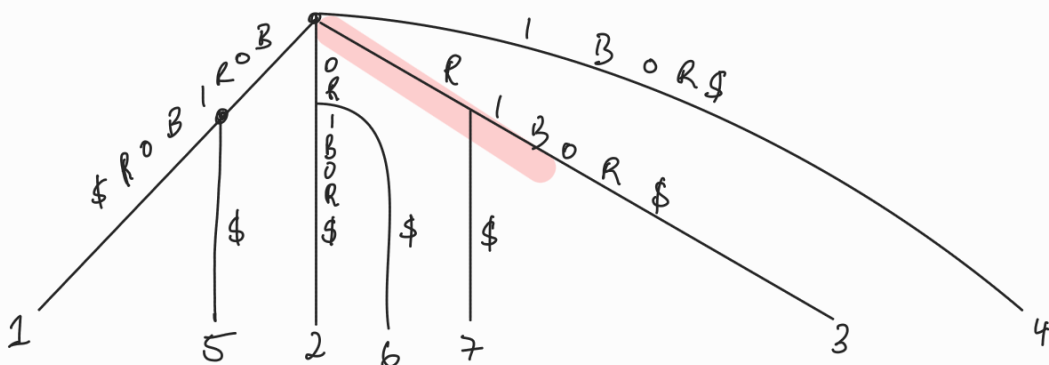
& $P[1..i]$ corresponds therefore to a prefix of some suffix of S , that is, a substring of S .

\Rightarrow if $i' = |P| \Rightarrow P$ occurs in T else, not.

\Rightarrow runtime $O(|P|)$

Exmpl: $T = \$ORIBOR\$$.

$P = RIB$



[every exist letter appears once at first letter of edge ($\$v$) for some v .]

Exmpl 2 "Substring - database - search"

Given strings $S_1 \dots S_e$ (Database of Texts)

Does P occur in at least one of the S_i ?

\Rightarrow let $\$1 \dots \e pairwise distinct symbols that do not occur in any S_i & P

put $S = S_1 \$1 S_2 \$2 \dots S_e \$e$

Construct suffix tree \mathcal{T} for S ($O(|S_1| + |S_2| + \dots + |S_e|)$ time)

use IDEA of Exmpl 1

Exmpl 3

Find all occurrences of P

given S with suffix tree \mathcal{T} & pattern P with $|P|=n$

Find all occurrences of P in \mathcal{T} , that is,

all indices i st $P = S[i:i+n-1]$

Call $FLP(\mathcal{T}, P)$

\rightarrow this gives $end = v$ or $end = \text{edge } (u, v)$

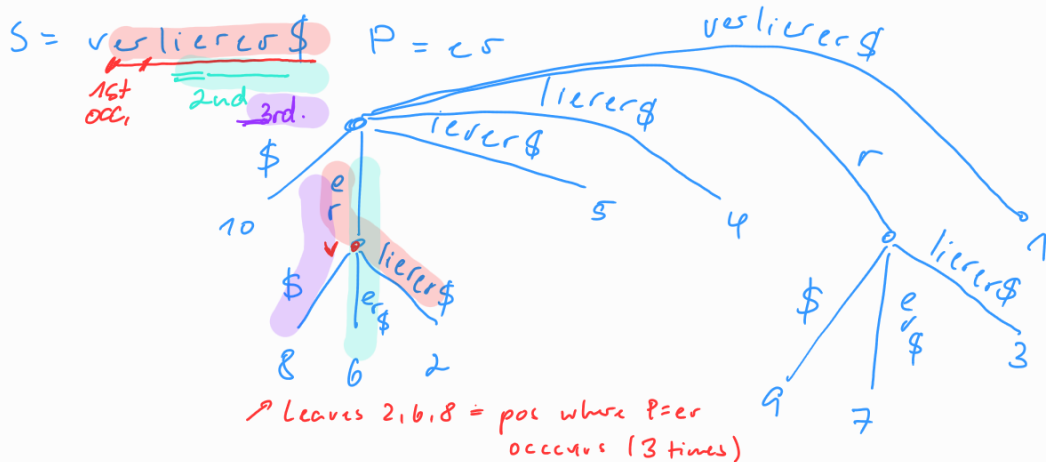


leaves $i \leq v = \#$ occurrences of P in S

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

[Exercise: works in $O(|P| + k)$ time]

$k = \text{nr of occurrences of } P$



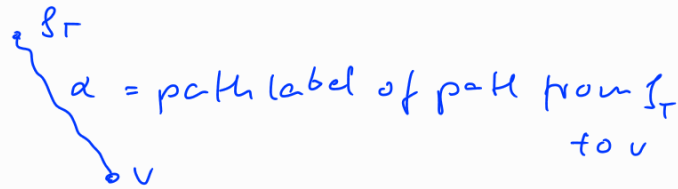
Exmpl 4

Find longest substring in S that occurs at least 2 times

[Application: Find repeated regions in genome]

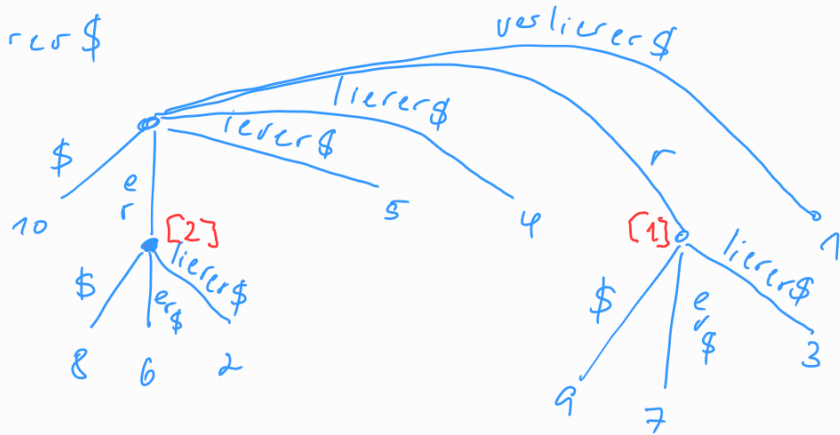
MISSISSIPPI issi

string depth of v in T is $|a|$



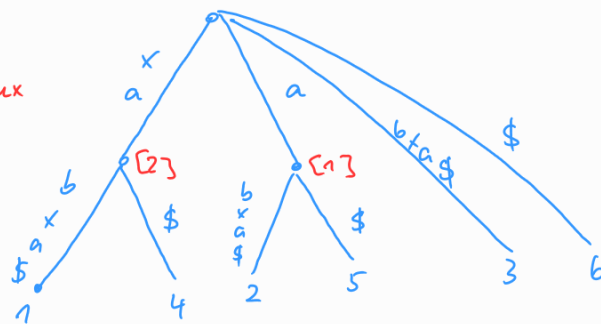
$S = \text{verlierer}\$$

string depth of inner vertex



$S = \text{xabxa}\$$

string depth of inner vertex



xa occurs 2 times.

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

\Rightarrow length of longest substring that occurs at least 2 times = maximum string-depth of inner vertices.

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

Exmpl 5

efficient detection of pairwise overlaps

Given $J = \{S_1 \dots S_N\}$ set of strings (eg. sequenced DNA fragments)

Aim: define $ov(S_i, S_j) \neq r_{ij}$
 "size of largest overlap"

$$S = S_1 \$1 S_2 \$2 \dots S_N \$N$$

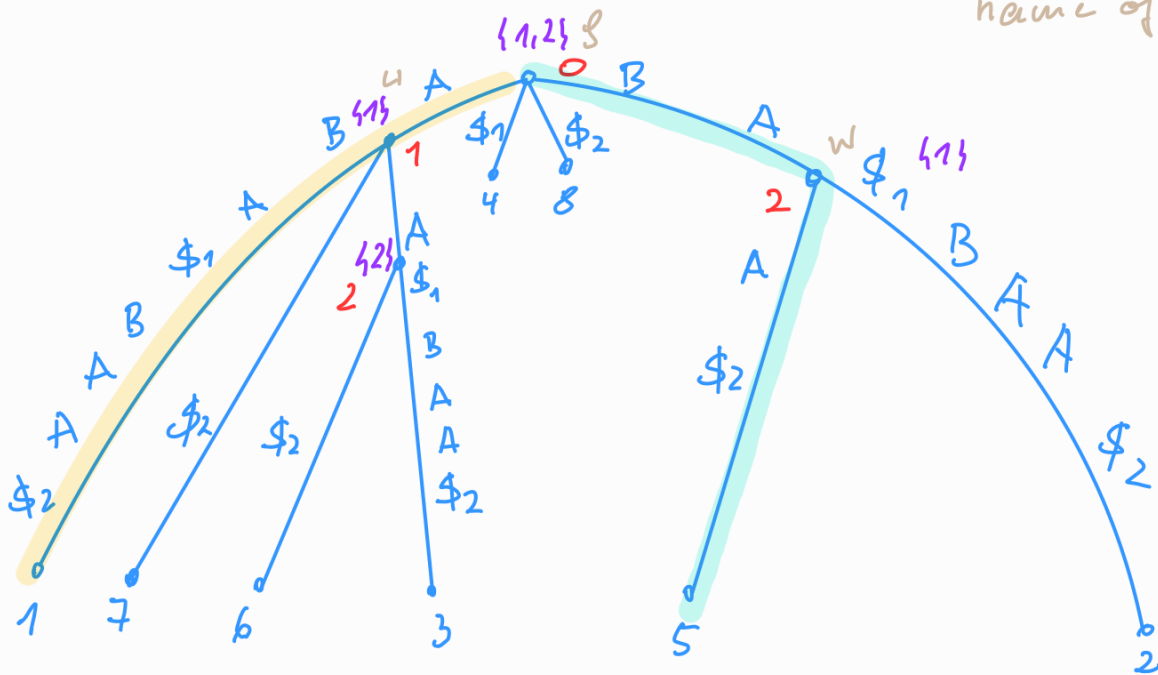


$$S_1 = ABA \rightarrow ov(S_1, S_2) = 2$$

$$S_2 = BAA \rightarrow ov(S_2, S_1) = 1$$

$$S = \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ A & B & A & \$1 & B & A & A & \$2 \end{matrix}$$

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 : $L(x) = \{i \mid \exists \text{ leaf } v \text{ st } (xv) \text{ edge in suffix tree} \& \text{ label } (xv) \text{ starts with } \$i\}$.

string depth $_{ot\ sr} =$ # characters on path from root to x .

For ($j=1 \dots N$)

$x = f(\text{root})$
WHILE (x not leaf)
 FOR (all $i \in L(x), i \neq j$)
 IF ($\text{depth}(x) < \min(|S_i|, |S_j|)$
 OR $(S_i, S_j) = \text{depth}(x)$)
 $x \leftarrow \text{child of } x \text{ on } \tau_j$

$j=1, x=f$

FOR ($i \in L(f) = \{1, 2\}, i \neq 1$)
 $\Rightarrow i=2$
 $\text{depth}(f) = 0 < \min(|S_1|, |S_2|)$
 $\Rightarrow \underline{\text{OR}(S_2, S_1) = 0}$

$x \leftarrow u$

FOR ($i \in L(u) = \{1\}, i \neq 1$)
 not any
 $x \leftarrow \text{leaf stop.}$

$j=2, x=f$

FOR ($i \in L(f) = \{1, 2\}, i \neq 2$)
 $\Rightarrow i=1$
 $\text{depth}(f) = 0 < \min(|S_1|, |S_2|)$
 $\text{OR}(S_1, S_2) = 0$

$x \leftarrow w$

FOR ($i \in L(w) = \{1\}, i \neq 2$)
 $\Rightarrow i=1$
 $\text{depth}(w) = 2 < \min(|S_1|, |S_2|)$
 $\Rightarrow \underline{\text{OR}(S_1, S_2) = 2}$
 $x \leftarrow \text{leaf stop.}$

[without proof of correctness or runtime analysis]

$O(N \|S\|)$

Exercise \rightarrow Book "Alg. Aspects of Bioning."

Appendix

Details of $O(n^2)$ Algorithm

SUFFIX-TREE(S)

init \mathcal{T} as  and label $(s_{\mathcal{T}}, 1) = S[1..m]$

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

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

IF (end = edge e) // e = (uw) with label γ

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



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

FIND-LONGEST-PATH (\mathcal{T} , S')

$j=1$, $v = s_{\mathcal{T}}$

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

Find edge $e = (vw)$ in \mathcal{T} , $w \in \mathcal{T}_v$, whose label starts with $S'(j)$

IF (such edge does not exist)

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

Let γ be label of e

$l = 1$

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

└ $j = j+1$
 $l = l+1$

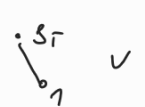
IF ($l \leq |\gamma|$) // while loop ended at some point before " $S'(j) = \gamma(l)$ "
 $\forall l = 1..|\gamma|$
 └ RETURN ($e, j-1, l-1$)
 ie, $S(j) \neq S(e)$

$v = w$ // go to consider edges starting at w

RETURN ($v, |S'|, \emptyset$)

Prop. 8.2

Alg. SUFFIXTREE correctly computes suffix tree for $S = x_1 \dots x_m$, $x_m = \$$.

proof: (X1): T_1 has 1 leaf, T_2 has 2 leaves ... T_m has m leaves ✓
 (sketch) (X2): T_1 : 

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



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

(X4) $T_i \checkmark$ Assume for $T_i \forall v$: label $(v_{i1}), (v_{i2})$ starts with different symbol, where v_1, v_2 are children of v .



Since α is longest substring that is prefix of

$$S[i+1..m] = \underbrace{x_{i+1} \dots x_r}_{\alpha} \dots \underbrace{x_{r+1} \dots x_m}_{\beta}$$

$\beta(1) \neq x_{r+1}$ as otherwise a not longest prefix! ✓

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

□

Runtime:

SUFFIXTREE(S), $S = x_1 \dots x_m$

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

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

→ m times FIND-LONGEST-PATH (FLP)
it called.

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

Find-edge → for v all neighbors w are considered

⇒ $\deg_{\rightarrow}(v)$ 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 ($j \in |S|$) DO

⋮

WHILE ($j \in |S|$ & $l \in |x|$ & $S'(j) = x(l)$) DO

$j = j + 1$
 ⋮

⇒ $j \leq m \Rightarrow O(m)$ time

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

⇒ FLP runs in $O(m)$ time

⇒ 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