

COL 351:

Analysis and Design of Algorithms

Lecture 14

String Matching

Given: String $S = [s_1, \dots, s_n]$
pattern $P = [p_1, \dots, p_k]$.

Find: Does there exists a
sub-string of S that is
identical to P.

Example:

$S =$ “cuckoo **hashing** is efficient”

$P =$ “hash”

Yes

Subproblem

Given: Pattern $P = [p_1, \dots, p_k]$.

Find: Table of size k satisfying

$Table[i] =$ Length of longest
common prefix and suffix
of $P[1, i]$

i	1	2	3	4	5	6
$P[i]$	A	B	C	A	B	B
$Table[i]$	0	0	0	1	2	0

Subproblem

Question:

Can we find **All** common prefixes/
suffixes of $P[1, i]$?



Eg. $P = (A B A B A B A B c)$

the length of ALL common
prefixes / suffixes of $P[1...8]$ is 6, 4, 2

Given: Pattern $P = [p_1, \dots, p_k]$.

Find: Table of size k satisfying

$Table[i] =$ Length of longest
common prefix and suffix
of $P[1, i]$

i	1	2	3	4	5	6
$P[i]$	A	B	C	A	B	B
$Table[i]$	0	0	0	1	2	0

Subproblem

Lemma: $\text{Table}[i]$, $\text{Table}[\text{Table}[i]]$, $\text{Table}[\text{Table}[\text{Table}[i]]]$,

gives the length of **ALL** common prefixes-suffixes of $P[1, i]$.

Proof idea:

$\text{Table}[18] = 8$

$P = (\overbrace{p_1 \ p_2 \ p_3 \ p_4 \ p_5 \ p_6 \ p_7 \ p_8}^Z \ p_9 \ p_{10} \ \overbrace{p_{11} \ p_{12} \ p_{13} \ p_{14} \ p_{15} \ p_{16} \ p_{17} \ p_{18}}^{i=18} \dots)$

Any common prefix/suffix of $P[1, 18]$

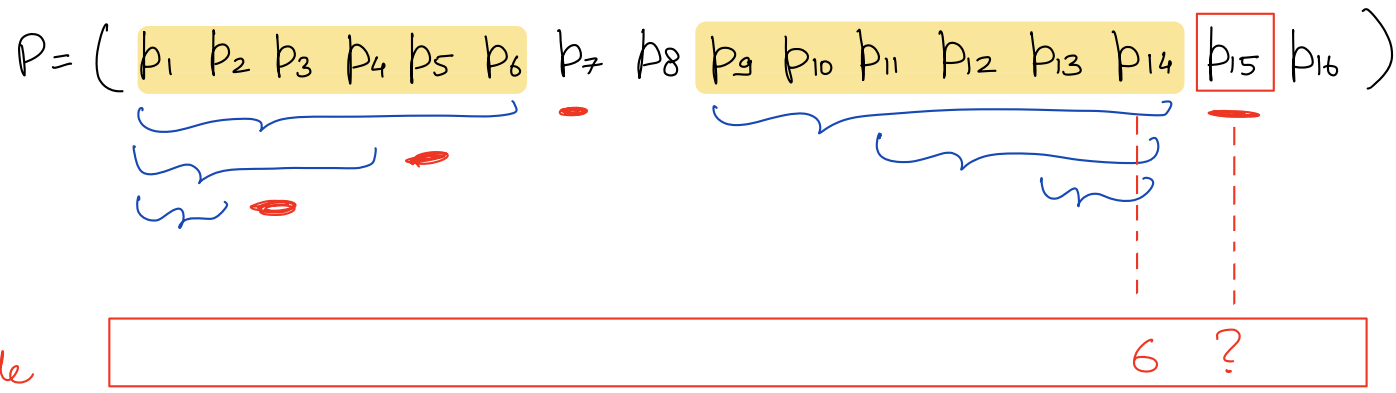
→ Either Z , or

→ A common prefix/suffix of Z

Subproblem

Lemma: $\text{Table}[i]$, $\text{Table}[\text{Table}[i]]$, $\text{Table}[\text{Table}[\text{Table}[i]]]$,
gives the length of **ALL** common prefixes-suffixes of $P[1, i]$.

Question: How to compute $\text{Table}[i + 1]$?



Subproblem

Lemma: $\text{Table}[i]$, $\text{Table}[\text{Table}[i]]$, $\text{Table}[\text{Table}[\text{Table}[i]]]$,

gives the length of **ALL** common prefixes-suffixes of $P[1, i]$.

Question: How to compute $\text{Table}[i + 1]$?

$L = \text{Table}[i];$

While ($L > 0$ and $P[i + 1] \neq P[L + 1]$): $L = \text{Table}[L];$

If ($P[i + 1] = P[L + 1]$): $L = L + 1;$

$\text{Table}[i + 1] = L;$

Values of L :

$\text{Table}[i], \text{Table}[\text{Table}[i]],$

$\text{Table}[\text{Table}[\text{Table}[i]]], \dots$

Stopping
Criteria

$L > 0$

$P[i + 1] = P[L + 1]$

$\text{Table}[i + 1] = L + 1$

$L = 0$

$P[i + 1] \neq P[L + 1]$

$\text{Table}[i + 1] = 1$

$L = 0$

$P[i + 1] \neq P[L + 1]$

$\text{Table}[i + 1] = 0$

Subproblem

$\text{Table}[i]$ = Length of longest (non-trivial) common prefix and suffix of $P[1, i]$

Table \leftarrow Array of size k ;

Table[1], $L = 0$;

For ($i = 1$ to $k - 1$):

/* value of L is Table[i] */

While ($L > 0$ and $P[i + 1] \neq P[L + 1]$): $L = \text{Table}[L]$;

If ($P[i + 1] = P[L + 1]$): $L = L + 1$;

Table[$i + 1$] = L ;

Subproblem

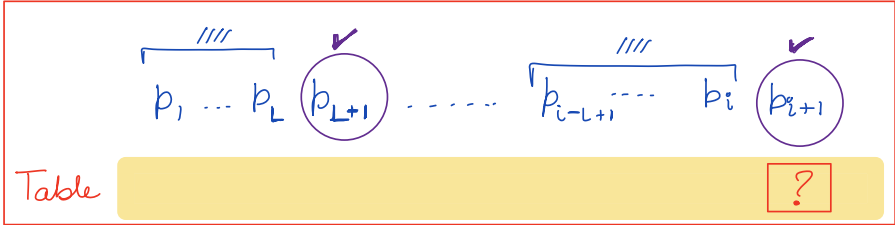
$\text{Table}[i]$ = Length of longest (non-trivial) common prefix and suffix of $P[1, i]$

```
Table ← Array of size k;  
Table[1], L = 0;  
For (i = 1 to k - 1):  
    /* value of L is Table[i] */  
    While (L > 0 and P[i + 1] ≠ P[L + 1]): L = Table[L];  
    If (P[i + 1] = P[L + 1]): L = L + 1;  
    Table[i + 1] = L;
```

Invariant:
When FOR LOOP starts
we have $L = \text{Table}[i]$.

Values of L :
 $\text{Table}[i], \text{Table}[\text{Table}[i]],$
 $\text{Table}[\text{Table}[\text{Table}[i]]], \dots$

STOPPING CRITERIA
 $P[i+1] = P[L+1]$
 $L=0 \begin{cases} P[1] = P[i+1] \\ \text{Table}[i+1] = 1 \\ P[1] \neq P[i+1] \\ \text{Table}[i+1] = 0 \\ L=0 \end{cases}$



Subproblem

$\text{Table}[i]$ = Length of longest (non-trivial) common prefix and suffix of $P[1, i]$

Table \leftarrow Array of size k ;

Table[1], $L = 0$;

For ($i = 1$ to $k - 1$):

/* value of L is Table[i] */

While ($L > 0$ and $P[i + 1] \neq P[L + 1]$): $L = \text{Table}[L]$;

If ($P[i + 1] = P[L + 1]$): $L = L + 1$;

Table[$i + 1$] = L ;

Analysis:

Claim 1: L is incremented at most k times.

Claim 2: Throughout the algorithm L can decrease at most k times, so total number of iterations of While loop is at most k .

Therefore, time complexity is $O(k)$.

Main Problem

Given: String $S = [s_1, \dots, s_n]$ and pattern $P = [p_1, \dots, p_k]$.

Find: Does there exists a **sub-string of S** that is identical to P.

Example:

$S =$ “cuckoo **hash**ing is efficient”

$P =$ “hash”

Yes

Main Problem

$\text{Table}[i] :=$ Length of longest (non-trivial) common prefix and suffix of $P[1, i]$

$A[i] :=$ Length of longest prefix of **P** that is also a suffix of **S** $[1, i]$

Eg:

$P =$ a b a b c

$S =$ s_1 s_2 s_3 s_4 a s_5 b s_6 a s_7 b s_8 s_9 s_{10}

⋮

A

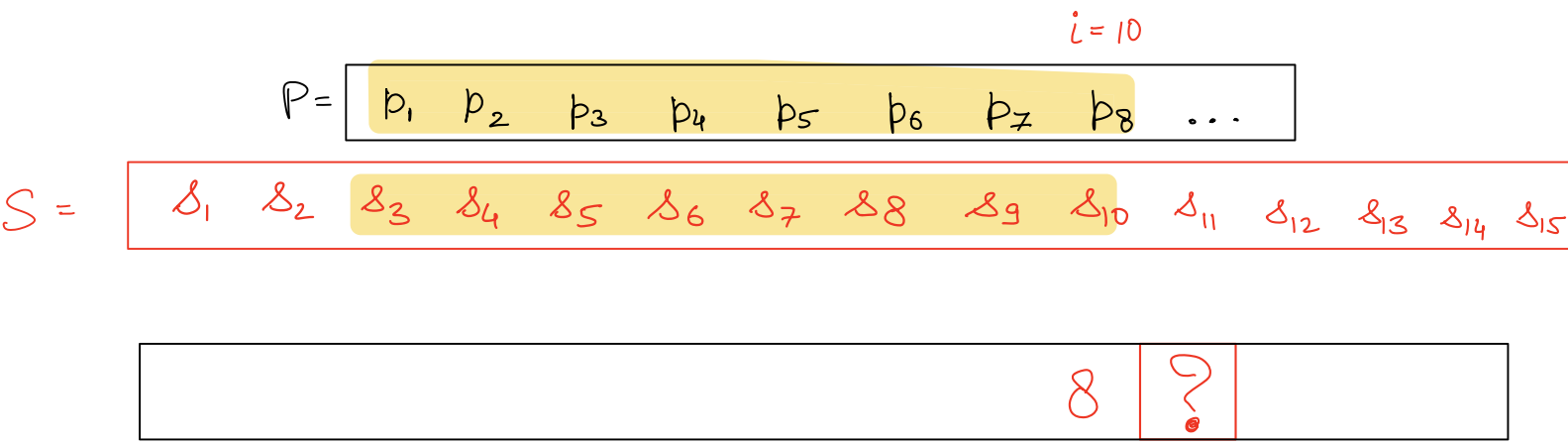
4

Main Problem

$\text{Table}[i] :=$ Length of longest (non-trivial) common prefix and suffix of $P[1, i]$

$A[i] :=$ Length of longest prefix of **P** that is also a suffix of **S** $[1, i]$

Question: How to find $A[i + 1]$?



Ques: How can you find All prefixes of P that are suffixes of $(s_1 \dots s_{10})$?

Ans: By sequence: $L = A[i], \text{Table}[L], \text{Table}[\text{Table}[L]], \dots$

Main Problem

$\text{Table}[i] :=$ Length of longest (non-trivial) common prefix and suffix of $P[1, i]$

$A[i] :=$ Length of longest prefix of **P** that is also a suffix of **S** $[1, i]$

***Question:** How to find $A[i + 1]$?*

$L = A[i];$

While ($L > 0$ and $S[i + 1] \neq P[L + 1]$): $L = \text{Table}[L];$

If ($S[i + 1] = P[L + 1]$): $L = L + 1;$

$A[i + 1] = L;$

Main Problem

$\text{Table}[i] :=$ Length of longest (non-trivial) common prefix and suffix of $P[1, i]$

$A[i] :=$ Length of longest prefix of **P** that is also a suffix of **S** $[1, i]$

$A \leftarrow$ Array of size $n + 1$ with $A[0] = 0$;

$L \leftarrow 0$;

For ($i = 0$ to $n - 1$):

 /* value of L is $A[i]$ */

 While ($L > 0$ and $\text{S}[i + 1] \neq \text{P}[L + 1]$): $\underline{L = \text{Table}[L]}$;

 If ($\text{S}[i + 1] = \text{P}[L + 1]$): $\underline{L = L + 1}$;

$A[i + 1] = L$;



**Knuth-Morris-Pratt
(KMP) algorithm**