



# SLIIT

*Discover Your Future*

# IT2070 – Data Structures and Algorithms

## Lecture 09

### String Matching Algorithms

### U. U. Samantha Rajapaksha

M.Sc.in IT, B.Sc.(Engineering) University of Moratuwa

Senior Lecturer SLIIT

[Samantha.r@slit.lk](mailto:Samantha.r@slit.lk)

• • • •

# Contents

- String Matching
- The naïve string matching algorithm
- The Rabin-Karp Algorithm
- Finite automata

we are given a pattern and we have to check if there is any given pattern

↓  
here we are going to check all these algorithms which help to check it

# String Matching

used to search a value in a string/text

- The problem of finding occurrence(s) of a pattern string within another string or body of text.
- There are many different algorithms for efficient searching.
- String matching is a very important subject in the wider domain of text processing.

# Applications

- In string matching problems, it is required to find the occurrences of a pattern in a text.
- **Applications**
  - Text processing
  - Text-editing e.g. *Find and Change* in word
  - Computer security (virus detection, password checking)
  - DNA sequence analysis.
  - Data communications (header analysis)



# Example

we keep this pattern on top of  
the text and check one by one

Text  $T$

A	B	C	A	B	A	A	B	C	A	B	A	C
---	---	---	---	---	---	---	---	---	---	---	---	---

Pattern  $P$

← Shift 0

A	B	A	A
---	---	---	---

A	B	A	A
---	---	---	---

A	B	A	A
---	---	---	---

Shift 1

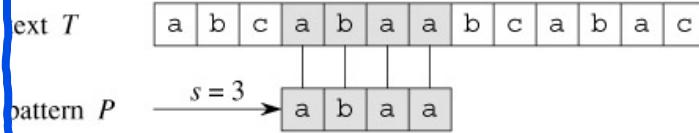
Shift 2

Shift 3

Checking starting from beginning  
shift by 1 index and keep in top of text

A	B	A	A
---	---	---	---

when shifted by 3 positions we find the fit



we have to continue this till the end of this arr.

# String Matching problem

- We formulate the ***String Matching problem*** as follows.
- We assume that the text is an array **T[1..n]** of length **n** and the pattern is an array **P[1..m]** of length **m**.
- We further assume that the elements of P and T are characters drawn from a finite alphabet  $\Sigma$ .
- For example we may have  $\Sigma=\{0,1\}$  or  $\Sigma=\{a,b,\dots,z\}$ .
- The character arrays P and T are often called ***Strings*** of characters.

# String Matching problem(contd.)

- We say that pattern  $P$  occurs with shift  $s$  in text  $T$  (or, equivalently that pattern  $P$  occurs beginning at position  $s + 1$  in text  $T$ )

number of shifts  
 $x \downarrow$   
 $0 \leq s \leq n-m$  and  $T[s+1..s+m] = P[1..m]$   
 length of pattern  
 length of text  
 within this range  
 If  $P$  occurs with shift  $s$  in  $T$ , then we call  $s$  a **valid shift**; otherwise we call  $s$  an **invalid shift**.

$T$ 

A	A	B	C	A	B	A
---	---	---	---	---	---	---

$P$ 

C	A	B
---	---	---

$n=7$   
 $n=3$

$s=0$ 

C	A	B
---	---	---

 X

$s=1$ 

C	A	B
---	---	---

 X

Pattern  $s=2$ 

C	A	B
---	---	---

 X

$s=3$ 

C	A	B
---	---	---

 ✓

because we can't go out of the arrays  $\rightarrow n-m$

$s=4$

C	A	B
---	---	---

 ✓

**Text  $T$**



**Pattern  $P$**

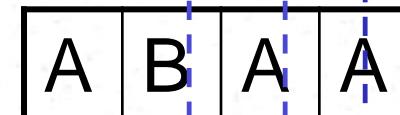
Shift 0



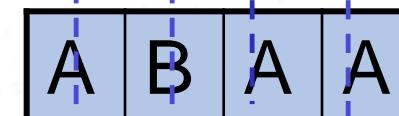
Shift 1



Shift 2



Shift 3



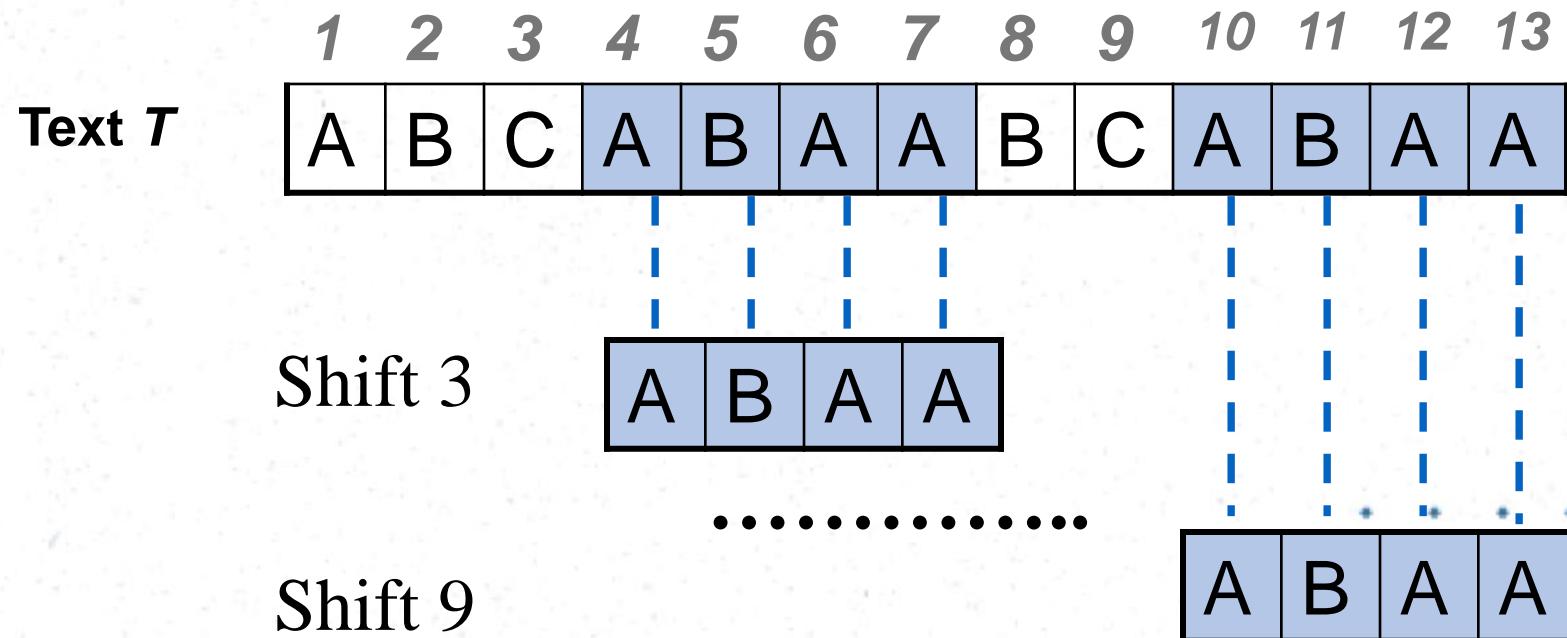
Here,  $n = 13$   $m = 4$   $s = 3$

$S = 3$  is a valid shift because

$0 \leq 3 \leq 13 - 4$  and  $T[3 + 1..3 + 4] = P[1..4]$ .

# String Matching problem(contd.)

The string-matching problem is the problem of finding all valid shifts with which a given pattern  $P$  occurs in a given text  $T$ .



# ① The naive string-matching algorithm

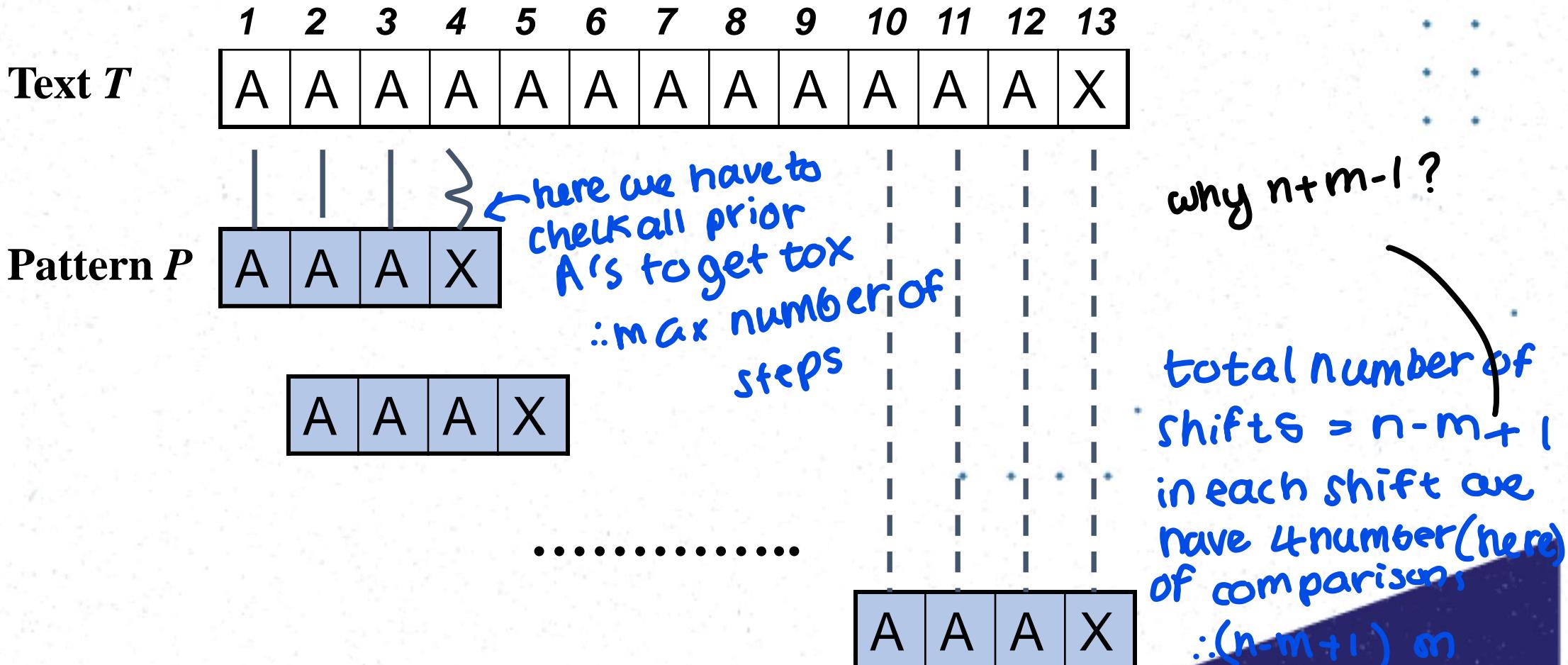
- The naïve algorithm finds all valid shifts using a loop that checks the condition  $P[1..m] = T[s+1..s+m]$  for each of the  $n-m+1$  possible values of  $s$ .

Naïve-String-Matcher ( $T, P$ )

```
1   $n = T.length$  length of text
2   $m = P.length$  length of pattern
3  for  $s = 0$  to  $\underline{n-m}$  from first shift to last shift
4    if  $P[1..m] = T[s+1..s+m]$  → here character by character we are checking
5      print "Pattern occurs with shift"  $s$ 
```

# When does worst case happen?

## Dworst case



# Worst case Analysis

## Running time

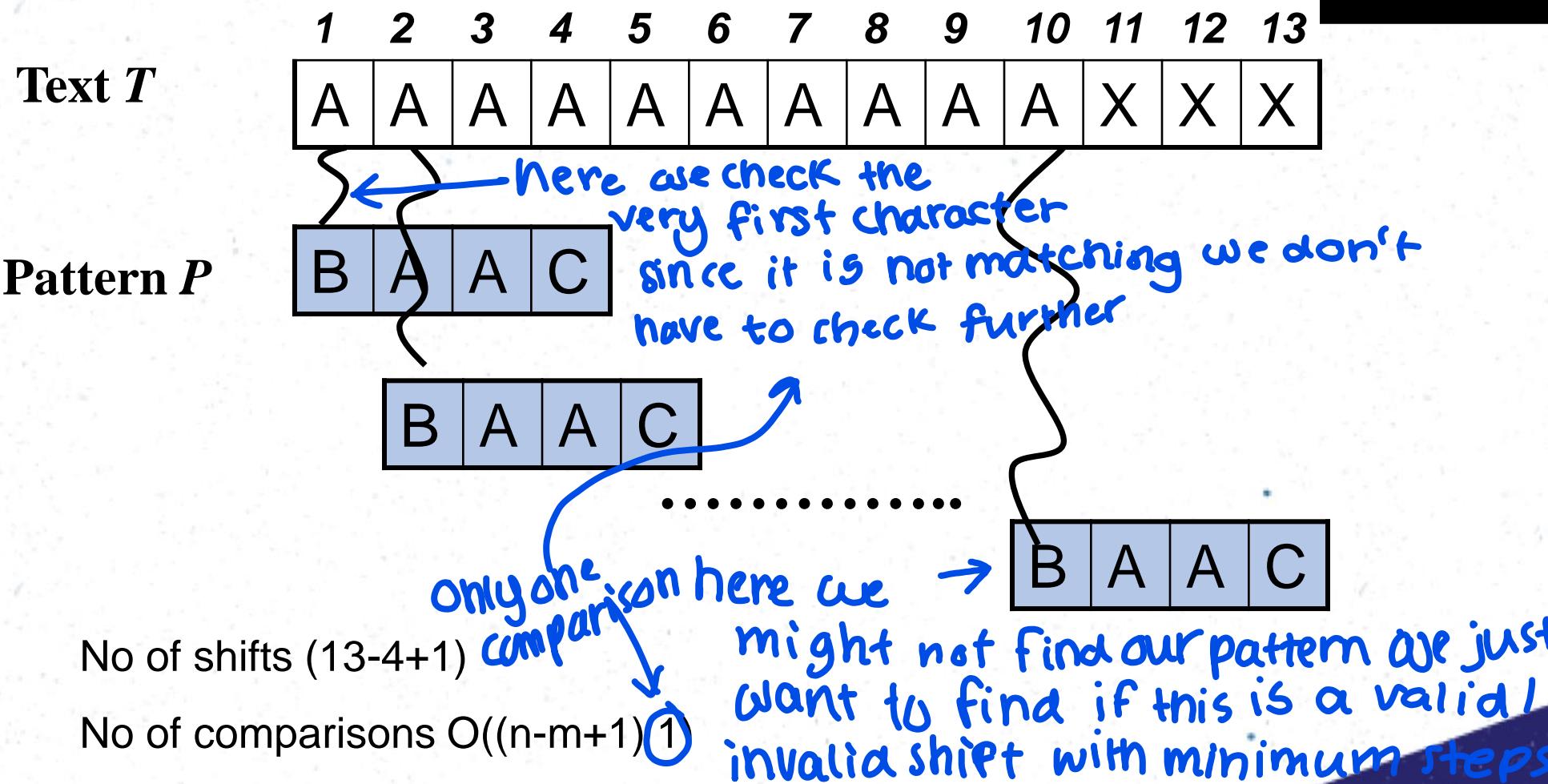
- This algorithm takes  $O((n-m+1)m)$  in the worst case
- There are at most  $n-m +1$  shifts
- $m$  is to compare each string after shifting
- Therefore worst case takes  $O((n-m+1)m)$
- E.g. In above example

$O( (13-4+1) 4 )$  comparisons

Shiftings

Comparing for each shift

## 2) Best case analysis



although the above algorithm is  
Simple but time consuming

# The Rabin-Karp Algorithm



**Professor.Richard Karp** Harvard University



**Professor Michel Rabin** Harvard University

Invented by Professor Richard Karp and Professor Michel Rabin in 1984.

① divide by modulo value and get the remainder

# The Rabin-Karp Algorithm

Given Text

2	3	5	9	0	2	3	1	4	1	5	2	6	7	3	9	9	2	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

$$13 \overline{)31415 \dots}$$

we get the  
remainder  
here it is 7

Find the occurrence of pattern

3	1	4	1	5
---	---	---	---	---

Using  
modulo 13

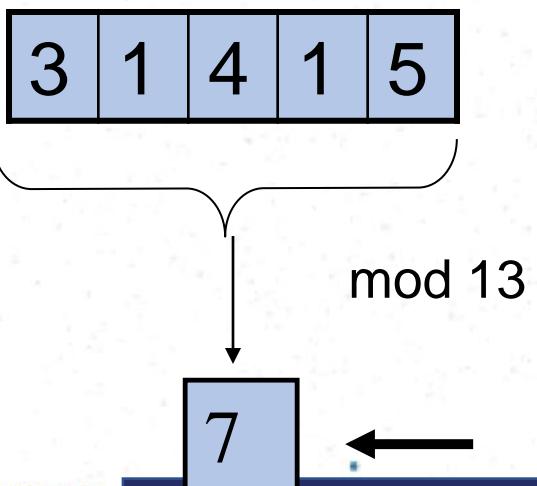
← we can decide  
this number  
modulo means the  
remainder

2	3	5	9	0	2	3	1	4	1	5	2	6	7	3	9	9	2	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

3	1	4	1	5
---	---	---	---	---

# Method

- Take the window size of the pattern
- Start from the beginning of the text taking windows
- Calculate the modulo value of that window
- Check it with the modulo value of the pattern

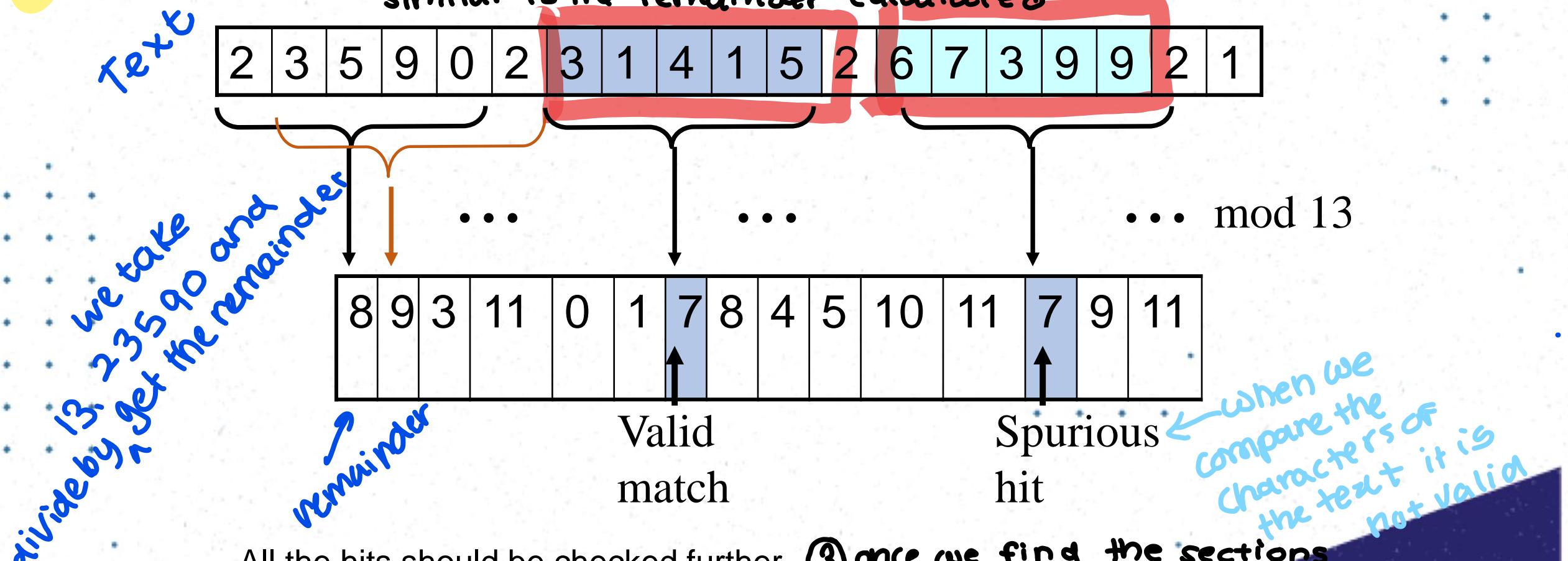


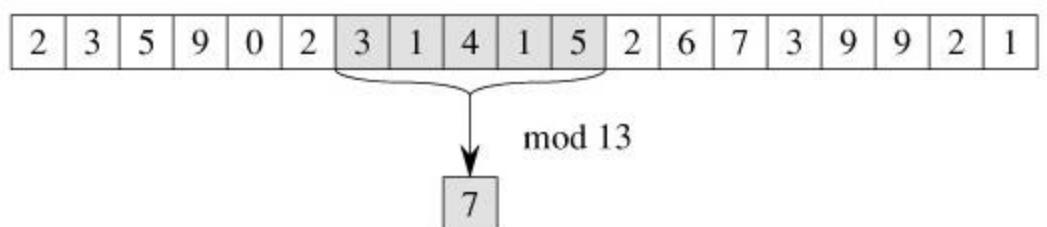
# The Rabin-Karp Algorithm(contd.)

- Let us assume that the input alphabet is  $\{0,1,2,\dots,9\}$ , so that each character is a decimal digit.
- We can then view a string of  $k$  consecutive characters as representing a length- $k$  decimal number.
- The character string 31425 thus corresponds to the decimal number 31,425.

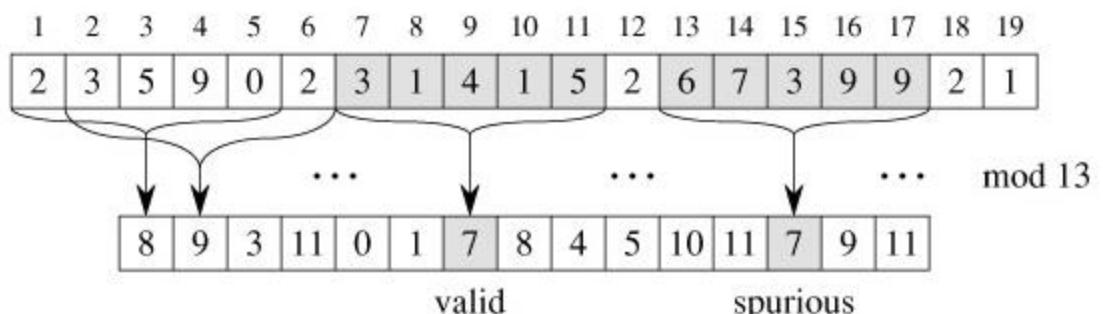
# The Rabin-Karp Algorithm

② repeatedly get the remainder and check if there are values similar to the remainder calculated

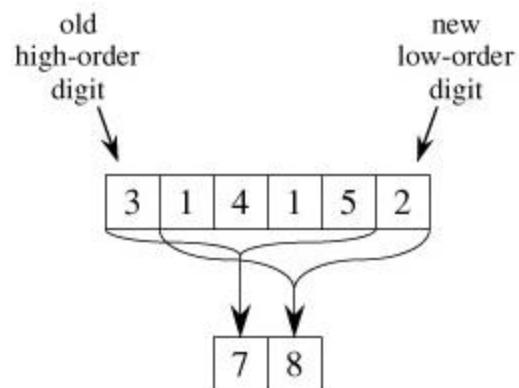




(a)



(b)



$$\begin{aligned}
 14152 &\equiv (31415 - 3 \cdot 10000) \cdot 10 + 2 \pmod{13} \\
 &\equiv (7 - 3 \cdot 3) \cdot 10 + 2 \pmod{13} \\
 &\equiv 8 \pmod{13}
 \end{aligned}$$

(c)

# Analysis of Rabin-Karp

- Comparing only the modulo value does not guarantee that the exact pattern is found.
- On the other hand, if modulo values are not matched, then we definitely know that it is not the pattern.
- All the hits must be tested further to see if the hit is a valid shift or just a spurious hit.
- This testing can be done by explicitly checking the condition  $P[1..m] = T[s + 1.. s + m]$
- If  $q$  is large enough, then we can hope that spurious hits occur infrequently enough that the cost of the extra checking is low.

*a is the modulus  
value*

*if a is small , you get somuch spurious hits*

# Worst case & Best case

## 1) worst case :

- If all the hits are spurious hits then we have to check each of those.
- Therefore the worst case occurs when all the hits are spurious hits

## 2) best case :

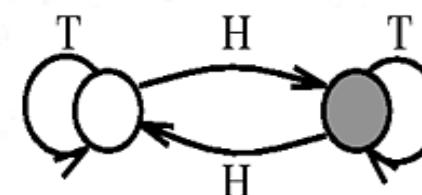
- If all the hits are neither spurious hits nor valid hits we don't have to check each of those.
- Therefore the best case occurs when no hits occurred

neither valid  
hitg / invalid hito

# String matching with finite automata

the way we find a given pattern in a string using finite automata

- Many string-matching algorithms build a finite automaton that scans the text string  $T$  for all occurrences of the pattern  $P$ .
- This section presents a method for building such an automaton.
- These string-matching automata are very efficient: they examine each text character *exactly once*, taking constant time per text character.



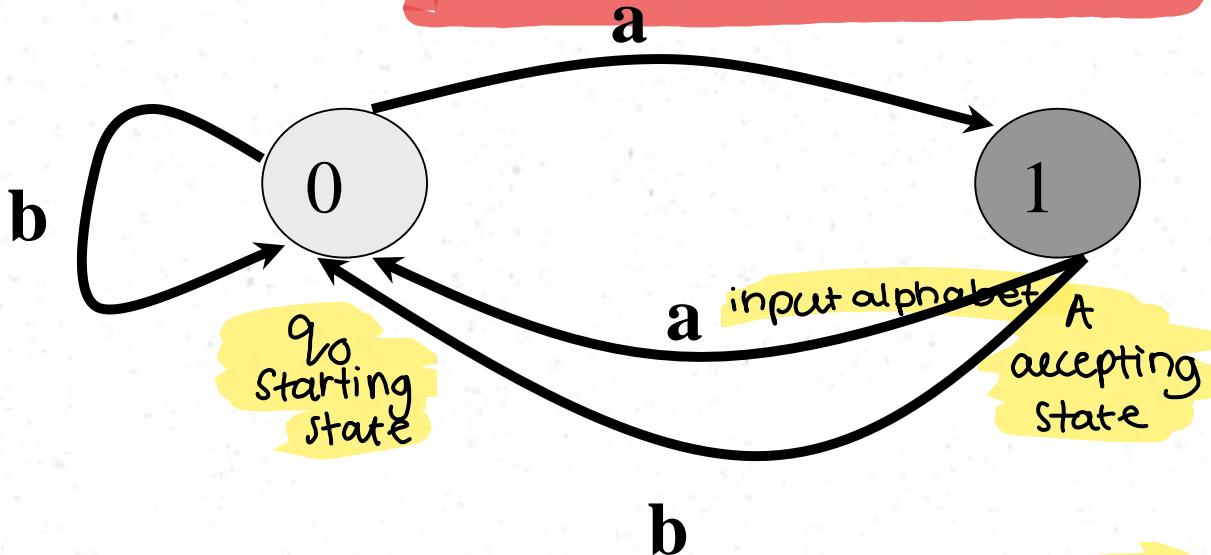
# Definition of a finite automaton

A ***finite automaton***  $M$  is a 5-tuple  $(Q, q_0, A, \Sigma, \delta)$ , where

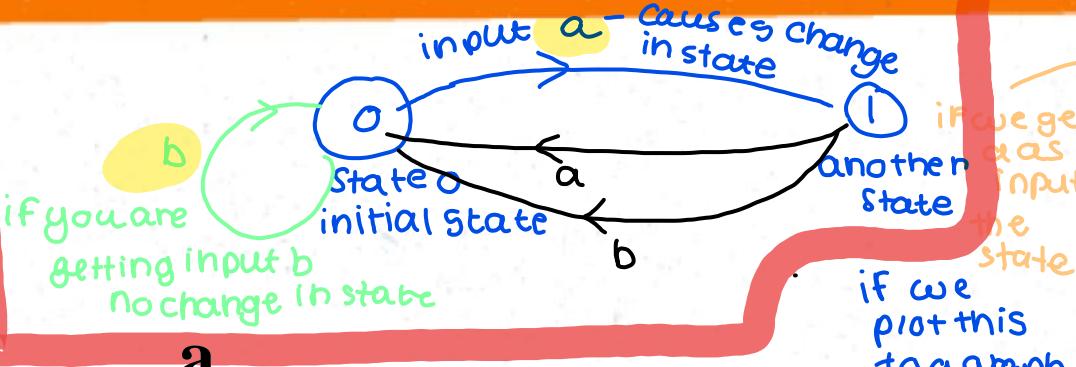
- $Q$  is a finite set of ***states***,
- $q_0 \in Q$  is the ***start state***,
- $A \subseteq Q$  is a distinguished set of ***accepting states***,
- $\Sigma$  is a finite ***input alphabet***,
- $\delta$  is a function from  $Q \times \Sigma$  into  $Q$ , called the ***transition function*** of  $M$ .



# Example



A simple two-state finite automaton with state set  $Q = \{0, 1\}$ , start state  $q_0 = 0$ , accepting state  $A = 1$ , input alphabet  $\Sigma = \{a, b\}$   $\rightarrow$  input value  
A tabular representation of the transition function  $\delta$

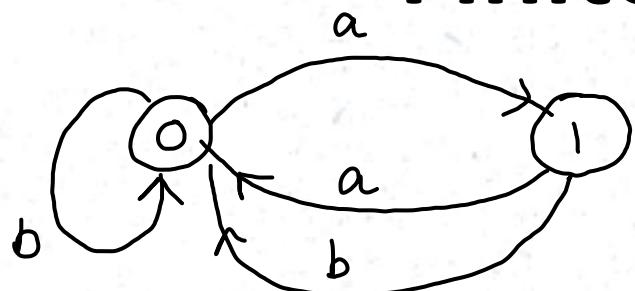


state	a	b
0	1	0
1	0	0

input

state	a	b
0	1	0
1	0	0

# Finite automata(contd.)



- Directed edges represent transitions.
- For example, the edge from state 1 to state 0 labeled **b** indicates  $\delta(1, b) = 0$ .
- This automaton accepts those strings that end in an odd number of **a**'s.
- For example, the sequence of states this automaton enters for input **abaaa** (including the start state) is  $(0, 1, 0, 1, 0, 1)$ , and so it accepts this input.
- For input **abbaa**, the sequence of states is  $(0, 1, 0, 0, 1, 0)$ , and so it rejects this input.

input 1 abaaa

assume that there is an input like  
this coming

input 2 abbaa

$0 \rightarrow 1 \rightarrow 0 \rightarrow 1 \rightarrow 0 \rightarrow 1$   
 $a \rightarrow b \rightarrow a \rightarrow a \rightarrow a$

$0 \rightarrow 1 \rightarrow 0 \rightarrow 0 \rightarrow 1 \rightarrow 0$   
 $a \rightarrow b \rightarrow b \rightarrow a \rightarrow a$

here since we move from state 0 to state 1 this has a change in state  
input is accepted  
since this doesn't have a change in state not valid  
 $\therefore$  input is rejected

here in this example we accept/rejects

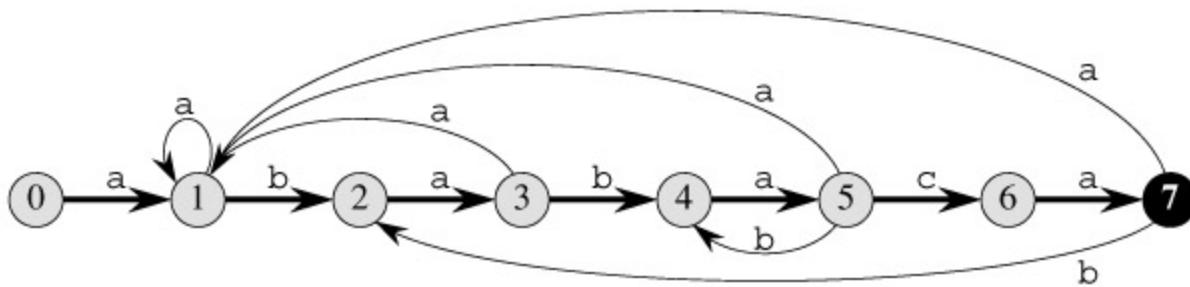
# Example. Accepts all strings ending in the string ababaca.

- (a) A state-transition diagram for the string-matching automaton that accepts all strings ending in the string ababaca. State 0 is the start state, and state 7 (shown blackened) is the only accepting state. A directed edge from state  $i$  to state  $j$  labeled  $a$  represents  $\delta(i, a) = j$ . The right-going edges forming the "spine" of the automaton, shown heavy in the figure, correspond to successful matches between pattern and input characters. The left-going edges correspond to failing matches. Some edges corresponding to failing matches are not shown; by convention, if a state  $i$  has no outgoing edge labeled  $a$  for some  $a \in \Sigma$ , then  $\delta(i, a) = 0$ .
- (b) The corresponding transition function  $\delta$ , and the pattern string  $P = \text{ababaca}$ . The entries corresponding to successful matches between pattern and input characters are shown shaded.
- (c) The operation of the automaton on the text  $T = \text{abababacaba}$ . Under each text character  $T[i]$  is given the state  $\varphi(T_i)$  the automaton is in after processing the prefix  $T_i$ . One occurrence of the pattern is found, ending in position 9.

Pattern = ababaca

Text = abababacaba

the states going to  $\emptyset$  not shown here



finite automata state transition diagram

state	input			$P$
	a	b	c	
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

(b)

Text	$i$	—	1	2	3	4	5	6	7	8	9	10	11
$T[i]$	—	a	b	a	b	a	b	a	c	a	b	a	
state $\phi(T_i)$	0	1	2	3	4	5	4	5	6	7	2	3	

(c)

# Summary

- String Matching
- The naïve string matching algorithm
- The Rabin-Karp Algorithm
- Finite automata