# String Matching

## String Matching

- T= text
- P=pattern
- n=text length
- m=pattern length.
- Z=alphabet size.

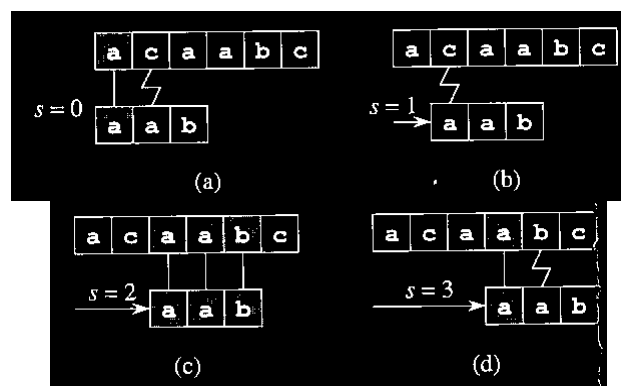# Naïve Method

3

---

# Naïve Method

$s = 0$ | a c a a b c / a a b

(a)

$s = 1$ | a c a a b c / a a b

(b)

$s = 2$ | a c a a b c / a a b

(c)

$s = 3$ | a c a a b c / a a b

(d)

**Complexity?**

# Knuth-Morris-Pratt Algorithm

5

## Idea of Jumping

(a)

(b)

# Jump Table

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $P[i]$ | a | b | a | b | a | b | a | b | c | a |
| $\pi[i]$ | 0 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 0 | 1 |

---

# KMP Algorithm

```
KMP-Matcher(T, P)
 1   n ← length[T]
 2   m ← length[P]
 3   π ← Compute-Prefix-Function(P)
 4   q ← 0
 5   for i ← 1 to n
 6       do while q > 0 and P[q + 1] ≠ T[i]
 7               do q ← π[q]
 8           if P[q + 1] = T[i]
 9               then q ← q + 1
10           if q = m
11               then print "Pattern occurs with shift" i − m
12                   q ← π[q]
```

## Computing Jumps

$$a \quad b \quad a \quad b \quad a \qquad P_q$$

$$a \quad b \quad a \qquad P_k$$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| $P[i]$ | a | b | a | b | a | b | a | b | c | a |
| $\pi[i]$ | 0 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 0 | 1 |

COMPUTE-PREFIX-FUNCTION($P$)
```
1   m ← length[P]
2   π[1] ← 0
    k ← 0
    for q ← 2 to m
        do while k > 0 and P[k + 1] ≠ P[q]
            do k ← π[k]
            if P[k + 1] = P[q]
                then k ← k + 1
            π[q] ← k
10  return π
```

- Jumps can be computed by preprocessing the pattern, by comparing the pattern with itself

---

## Complexity

- The running time for COMPUTE-PREFIX-FUNCTION is:
  - O(m)
- The total jump cannot exceed (m+n),= Thus complexity is O(m+n).

# State Machine

---
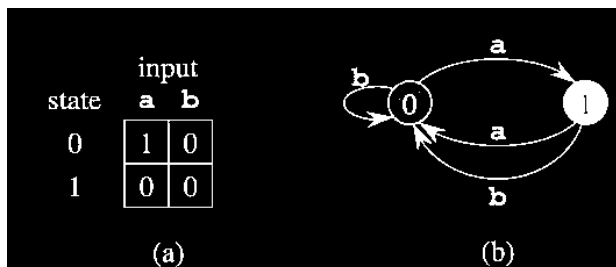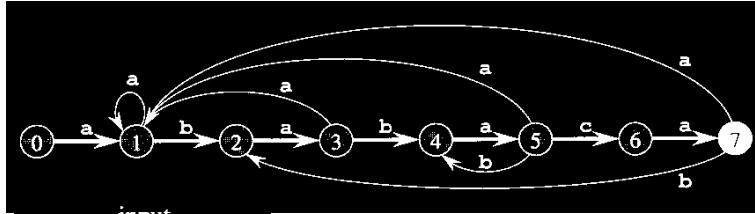
## FSM for Detecting String Which Ends with Even 1s

**Figure 34.5** A simple two-state finite automaton with state set $Q = \{0, 1\}$, start state $q_0 = 0$, and input alphabet $\Sigma = \{a, b\}$. (a) A tabular representation of the transition function $\delta$. (b) An equivalent state-transition diagram. State 1 is the only accepting state (shown blackened). 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. More precisely, a string $x$ is accepted if and only if $x = yz$, where $y = \varepsilon$ or $y$ ends with a b, and $z = a^k$, where $k$ is odd. For example, the sequence of states this automaton enters for input abaaa (including the start state) is $\langle 0, 1, 0, 1, 0, 1 \rangle$, and so it accepts this input. For input abbaa, the sequence of states is $\langle 0, 1, 0, 0, 1, 0 \rangle$, and so it rejects this input.

State Transition Diagram for String Matching

| state | input a | b | c | P |
|-------|---------|---|---|---|
| 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 | |

| $i$ | — | 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 |

*Operation of FSM on string abababacaba*

---

# Complexity

- Fastest Compute State Transition is:
  - O(m z)
- Total WC running time is:
  - O(n+mz)

# Boyer-Moore Algorithm
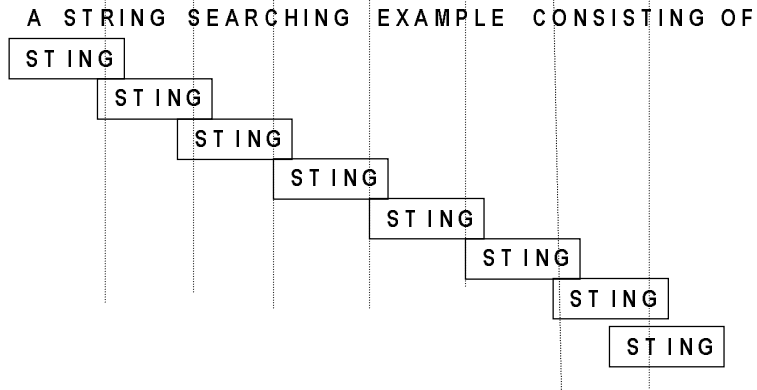# (1976)

## Boyer-Moore Algorithm

- Comparing from the Right to Left:
  - in the pattern, each time, there is a mismatch, see how many position the pattern can be shifted left.

- More Look ahead in the Preprocessing
  - bring into consideration the character that caused the mismatch while considering what to do next.

# Example

A STRING SEARCHING EXAMPLE CONSISTING OF

```
ST ING
    ST ING
        ST ING
            ST ING
                ST ING
                    ST ING
                        ST ING
                            ST ING
```

DESIGN &
ANALYSIS OF
ALGORITHM

# Complexity

- Boyer-Moore string search algorithm never uses more than M+N character comparisons, and uses about N/M steps of the alphabet is not small and the pattern is not long.

DESIGN &
ANALYSIS OF
ALGORITHM

# Rabin-Karp Method (1980)

## Rabin-Karp Algorithm

# Computing Hash Value



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

---

# RK Algorithm

```
RABIN-KARP-MATCHER(T, P, d, q)
1    n ← length[T]
2    m ← length[P]
3    h ← d^{m-1} mod q
4    p ← 0
5    t_0 ← 0
6    for i ← 1 to m
7        do p ← (dp + P[i]) mod q
8           t_0 ← (dt_0 + T[i]) mod q
9    for s ← 0 to n − m
10       do if p = t_s
11          then if P[1..m] = T[s + 1..s + m]
12               then "Pattern occurs with shift" s
13          if s < n − m
14             then t_{s+1} ← (d(t_s − T[s + 1]h) + T[s + m + 1]) mod q
```

- Radix is d. The prime is q.

# Complexity

- In the worst case the running time is $O((n-m+1)\,m)$.
    - (case $T=a^n$ and $P=a^m$)
    - Each evaluation after the first one is $O(1)$ in text.
- Average Case Complexity
    - Only one match in most cases $O(1)$
    - Thus running time is $O(n+m)$.