

Analysis of Algorithms

String Matching

1

3/3/2003

The Knuth-Morris-Pratt Algorithm

- String matching algorithm that runs in linear time ($O(n+m)$) by avoiding the computation of the transition function δ
 - Pattern matching is done using an auxiliary prefix function $\pi[1..m]$ precomputed from the pattern in time $O(m)$.
 - The array $\pi[1..m]$ allows an efficient computation “on-the-fly” of the transition function δ .

2

3/3/2003

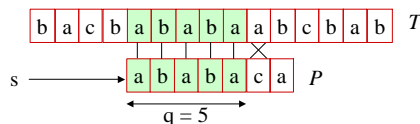
The Knuth-Morris-Pratt Algorithm

- For any state $q = 0, 1, \dots, m$, and any character $c \in \Sigma$, the value $\pi[q]$ contains the information needed to compute $\delta(q, c)$ that is independent of c .
- Prefix function $\pi[q] \Rightarrow O(m)$ (substantial savings, particularly if $|\Sigma|$ is large)
- Transition function $\delta[q, c] \Rightarrow O(m|\Sigma|)$

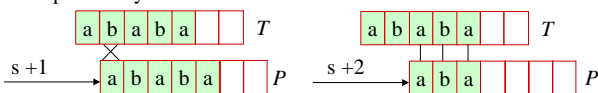
3

3/3/2003

Key Idea



- Using only our knowledge that **the 5 first characters matched**, we can deduce that a shift $s+1$ is invalid and that a shift $s+2$ is potentially valid.



Thus, we can safely transition to state $q=3$ with a shift $s' = s+2$

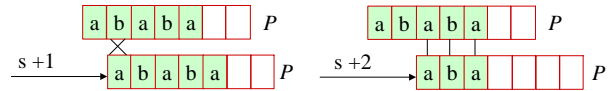
- $\pi[5] = 3$ and $s' = s + (q - \pi[q])$

4

3/3/2003

The Prefix Function for a Pattern

- Function π can be computed by comparing the pattern against itself



$$\Rightarrow \pi[5] = 3 \text{ and } s' = s + (q - \pi[q])$$

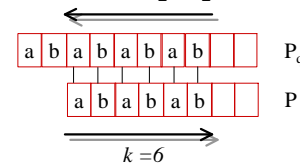
5

3/3/2003

The Prefix Function for a Pattern

- Formally, π is a function $\{1, 2, \dots, m\} \rightarrow \{0, 1, \dots, m-1\}$ such that

$$\pi[q] = \max \{k : k < q \text{ and } P_k \supset P_q\}$$
- That is, $\pi[q]$ is the length k of the longest prefix of P that is a **proper** suffix of P_q



6

3/3/2003

Prefix Function Computation

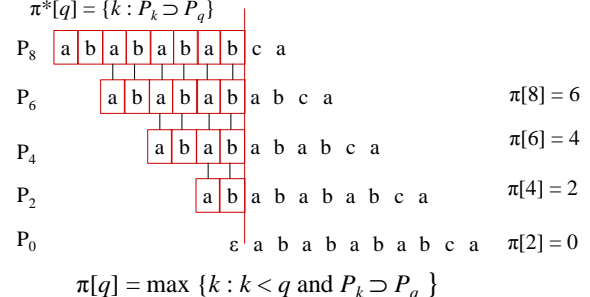
- We will show that by iterating the prefix function π , we can enumerate all the prefixes P_k that are suffixes of a given prefix P_q .
- Let $\pi^*[q] = \{q, \pi[q], \pi^2[q], \dots, \pi^l[q]\}$
- where $\pi^l[q]$ is defined in terms of functional composition, so that
 - $\pi^0[q] = q$
 - $\pi^{i+1}[q] = \pi[\pi^i[q]]$ for $i > 0$
- and the sequence in $\pi^*[q]$ stops when $\pi^l[q] = 0$.

7

3/3/2003

Enumerating all the prefixes P_k that are suffixes of a prefix P_q via π^*

- Prefix-function iteration lemma:** Let P be a pattern of length m with a prefix function π . Then, for $q = 1, 2, \dots, m$, we have $\pi^*[q] = \{k : P_k \supset P_q\}$



$$\pi[q] = \max \{k : k < q \text{ and } P_k \supset P_q\}$$

8

3/3/2003

Use of π^* to Compute the Prefix Function π

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

9

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$

$q=2$ ($k=0$)

a b a b a b a b c a

$P[k+1] \neq P[q]$

$k=0 \Rightarrow$ no smaller prefixes.

$\pi[q] \leftarrow k \Rightarrow \pi[2] \leftarrow 1$

10

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$
 $\pi[2] \leftarrow 0$

$q=3$ ($k=0$)

a b a b a b a b c a

$P[k+1] = P[q]$

$k = k + 1 = 1$ ($P_k \Rightarrow P_1$) .

$\pi[q] \leftarrow k \Rightarrow \pi[3] \leftarrow 1$

P_1 is the longest prefix that is proper suffix of P_3

11

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$
 $\pi[2] \leftarrow 0$
 $\pi[3] \leftarrow 1$

$q=4$ ($k=1$)

a b a b a b a b c a

$P[k+1] = P[q]$

$k = k + 1 = 2$ ($P_k \Rightarrow P_2$) .

$\pi[q] \leftarrow k \Rightarrow \pi[4] \leftarrow 2$

P_2 is the longest prefix that is proper suffix of P_4

12

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$
 $\pi[2] \leftarrow 0$
 $\pi[3] \leftarrow 1$
 $\pi[4] \leftarrow 2$

$q=5$ ($k=2$)

a b a b a b a b c a

$P[k+1] = P[q]$

$k = k + 1 = 3$ ($P_k \Rightarrow P_3$) .

$\pi[q] \leftarrow k \Rightarrow \pi[5] \leftarrow 3$

P_3 is the longest prefix that is proper suffix of P_5

a b a b a b a b c a

13

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$
 $\pi[2] \leftarrow 0$
 $\pi[3] \leftarrow 1$
 $\pi[4] \leftarrow 2$
 $\pi[5] \leftarrow 3$

$q=6$ ($k=3$)

a b a b a b a b c a

$P[k+1] = P[q]$

$k = k + 1 = 4$ ($P_k \Rightarrow P_4$) .

$\pi[q] \leftarrow k \Rightarrow \pi[6] \leftarrow 4$

P_4 is the longest prefix that is proper suffix of P_6

a b a b a b a b c a

14

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$ $\pi[6] \leftarrow 4$
 $\pi[2] \leftarrow 0$
 $\pi[3] \leftarrow 1$
 $\pi[4] \leftarrow 2$
 $\pi[5] \leftarrow 3$

$q=7$ ($k=4$)

a b a b a b a b c a

$P[k+1] = P[q]$

$k = k + 1 = 5$ ($P_k \Rightarrow P_5$) .

$\pi[q] \leftarrow k \Rightarrow \pi[7] \leftarrow 5$

P_5 is the longest prefix that is proper suffix of P_7

a b a b a b a b c a

15

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$ $\pi[6] \leftarrow 4$
 $\pi[2] \leftarrow 0$ $\pi[7] \leftarrow 5$
 $\pi[3] \leftarrow 1$
 $\pi[4] \leftarrow 2$
 $\pi[5] \leftarrow 3$

$q=8$ ($k=5$)

a b a b a b a b c a

$P[k+1] = P[q]$

$k = k + 1 = 6$ ($P_k \Rightarrow P_6$)

$\pi[q] \leftarrow k \Rightarrow \pi[8] \leftarrow 6$

P_6 is the longest prefix that is proper suffix of P_8

a b a b a b a b c a

16

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$ $\pi[6] \leftarrow 4$
 $\pi[2] \leftarrow 0$ $\pi[7] \leftarrow 5$
 $\pi[3] \leftarrow 1$ $\pi[8] \leftarrow 6$
 $\pi[4] \leftarrow 2$
 $\pi[5] \leftarrow 3$

$q = 9$ ($k = 6$)

a b a b a b a b c a

$P[k+1] \neq P[q]$

a b a b a b a b c
a b a b a b a

17

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$ $\pi[6] \leftarrow 4$
 $\pi[2] \leftarrow 0$ $\pi[7] \leftarrow 5$
 $\pi[3] \leftarrow 1$ $\pi[8] \leftarrow 6$
 $\pi[4] \leftarrow 2$
 $\pi[5] \leftarrow 3$

$q = 9$ ($k = 4$)

a b a b a b a b c a

$P[k+1] \neq P[q]$

a b a b a b a b c
a b a b a b a

18

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$ $\pi[6] \leftarrow 4$
 $\pi[2] \leftarrow 0$ $\pi[7] \leftarrow 5$
 $\pi[3] \leftarrow 1$ $\pi[8] \leftarrow 6$
 $\pi[4] \leftarrow 2$
 $\pi[5] \leftarrow 3$

$q = 9$ ($k = 2$)

a b a b a b a b c a

$P[k+1] \neq P[q]$

a b a b a b a b c
a b a

19

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$ $\pi[6] \leftarrow 4$
 $\pi[2] \leftarrow 0$ $\pi[7] \leftarrow 5$
 $\pi[3] \leftarrow 1$ $\pi[8] \leftarrow 6$
 $\pi[4] \leftarrow 2$
 $\pi[5] \leftarrow 3$

$q = 9$ ($k = 0$)

a b a b a b a b c a

$P[k+1] \neq P[q]$

a b a b a b a b c
a

20

3/3/2003

Compute-Prefix-Function(P)

```

 $m \leftarrow \text{length}[P]$ 
 $\pi[1] \leftarrow 0$   $\triangleright$  true for any pattern
 $k \leftarrow 0$ 
for  $q \leftarrow 2$  to  $m$ 
  do while  $k > 0$  and  $P[k+1] \neq P[q]$ 
    do  $k \leftarrow \pi[k]$ 
  if  $P[k+1] = P[q]$ 
    then  $k = k + 1$ 
   $\pi[q] \leftarrow k$ 
return  $\pi$ 

```

$\pi[1] \leftarrow 0$ $\pi[6] \leftarrow 4$
 $\pi[2] \leftarrow 0$ $\pi[7] \leftarrow 5$
 $\pi[3] \leftarrow 1$ $\pi[8] \leftarrow 6$
 $\pi[4] \leftarrow 2$ $\pi[9] \leftarrow 0$
 $\pi[5] \leftarrow 3$

$q = 10$ ($k = 0$)

a b a b a b a b c a

$P[k+1] = P[q]$

a b a b a b a b c a
a

$k = k + 1 = 1$ ($P_k \Rightarrow P_1$)

$\pi[q] \leftarrow k \Rightarrow \pi[10] \leftarrow 1$

21

3/3/2003

The KMP Algorithm

KMP-Matcher(T, P)

```

 $n \leftarrow \text{length}[T]$ 
 $m \leftarrow \text{length}[P]$ 
 $\pi \leftarrow \text{Compute-Prefix-Function}(P)$ 
 $q \leftarrow 0$   $\triangleright$  current position in  $P$ 
for  $i \leftarrow 1$  to  $n$   $\triangleright$  current position in  $T$ 
  do while  $q > 0$  and  $P[q+1] \neq T[i]$ 
    do  $q \leftarrow \pi[q]$ 
  if  $P[q+1] = T[i]$ 
    then  $q = q + 1$ 
  if  $q = m$ 
    then print "Pattern occurs with shift"  $i - m$ 
     $q \leftarrow \pi[q]$ 

```

22

3/3/2003