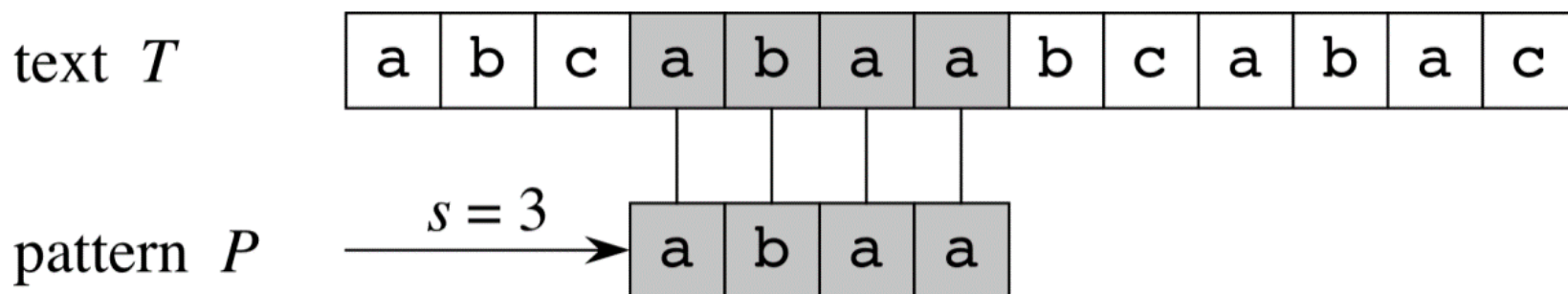


STRING MATCHING

Prof. Michael Tsai

2025/04

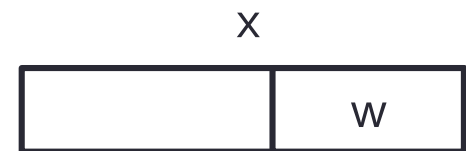
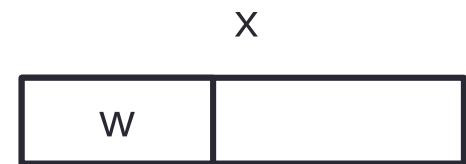
問題定義：字串比對



- 陣列 $T[1:n]$ 中有一個長度為 n 的字串
- 陣列 $P[1:m]$ 中有一個長度為 m 的字串 $m \leq n$
- P 和 T 的字串從一個字元的集合 Σ 中拿出
($\Sigma = \{0,1\}$ 或 $\Sigma = \{a, b, \dots, z\}$)
- 在 T 中間找到**所有** P 出現的位置 (valid shifts)
- Pattern P **occurs with shift s** in text T (beginning at position $s+1$)
 \rightarrow if $T[s + j] == P[j]$, for $1 \leq j \leq m$.
- If P occurs with shift s in T , we call s a **valid shift**.
 Otherwise, we call s an **invalid shift**.

一些定義

- Σ^* : 所有使用 Σ 中字元組成的有限長度字串 (包括長度為0的空字串)
- $|x|$: 字串 x 的長度
- xy : 把字串 x 和 y 接起來 (concatenation)
- $w \sqsubset x$: 字串 w 是字串 x 的prefix (也就是 $x=wy, y \in \Sigma^*$)
($w \sqsubset x$ 表示 $|w| \leq |x|$)
- $w \sqsupset x$: 字串 w 是字串 x 的suffix (也就是 $x=yw, y \in \Sigma^*$)
($w \sqsupset x$ 表示 $|w| \leq |x|$)
- 例: $ab \sqsubset abcca, cca \sqsupset abcca$



一些定義

- 空字串 ϵ 為任何字串的prefix & suffix
- 對任何字串 x, y 和字元 a , $x \sqsubset y$ iff $xa \sqsubset ya$
- \sqsubset 和 \sqsupset 為transitive(具遞移律)的operator,
e.g., $x \sqsubset y$ and $y \sqsubset z \Rightarrow x \sqsubset z$

iff means "if and only if"

Running Time of Different String-Matching Algorithms

Algorithm	Preprocessing time	Matching time
Naïve	$O(1)$	$O((n - m + 1)m)$
Rabin-Karp	$\Theta(m)$	$O((n - m + 1)m)$
Knuth-Morris-Pratt	$\Theta(m)$	$O(n)$

方法一：暴力法 (Naïve)

Cormen 32.1

Native-String-Matcher(T, P, n, m)

```
/* n=T.length
```

```
   m=P.length */
```

```
for s=0 to n-m
```

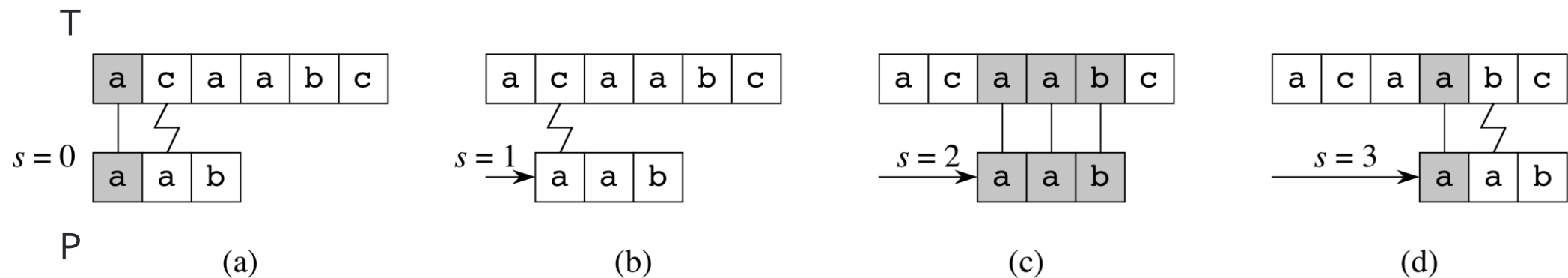
$n-m+1$ 次

```
    if P[1:m]==T[s+1:s+m]
```

```
        print "Pattern occurs with shift" s
```

$O(m)$

$O((n - m + 1)m)$

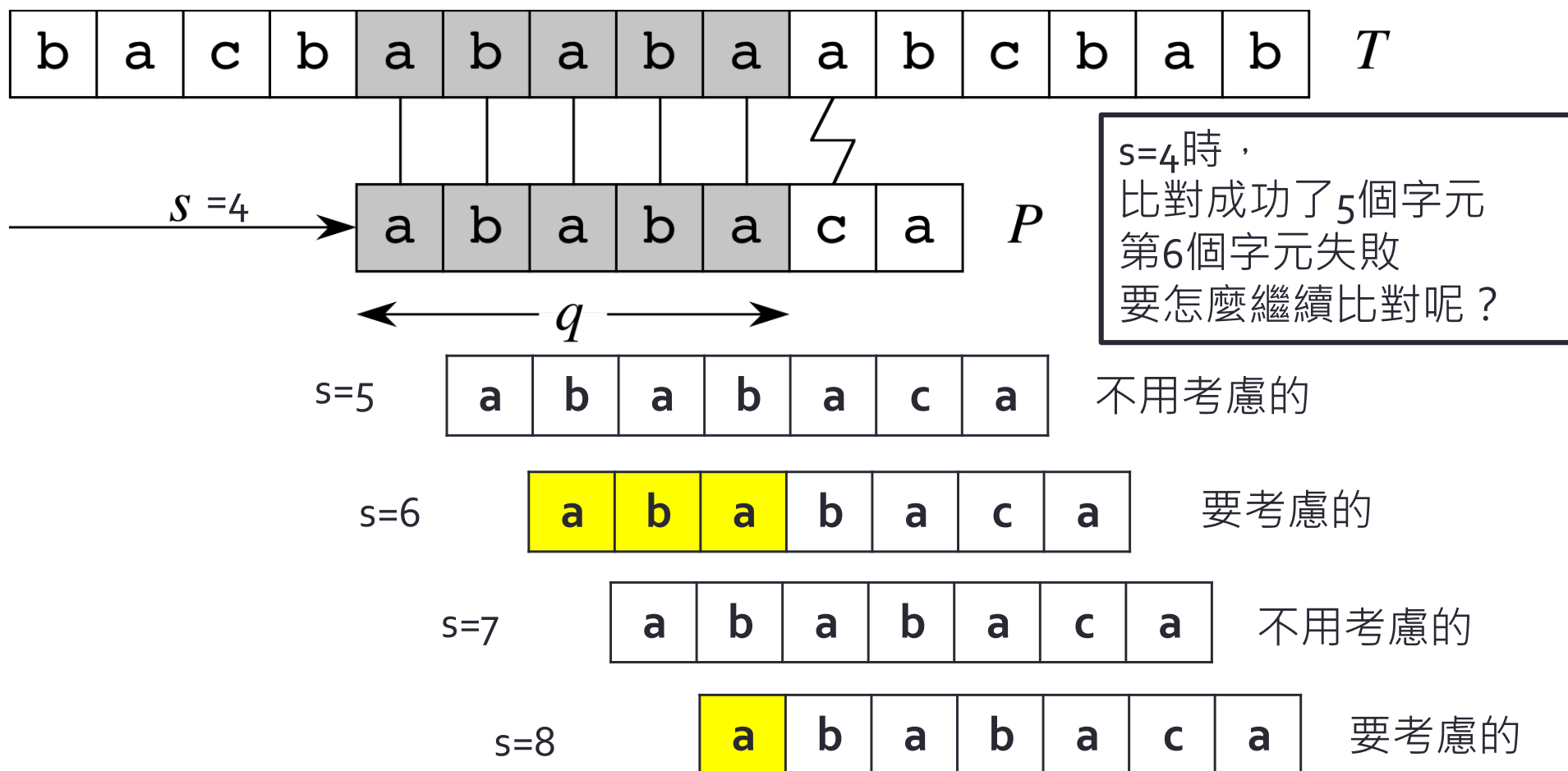


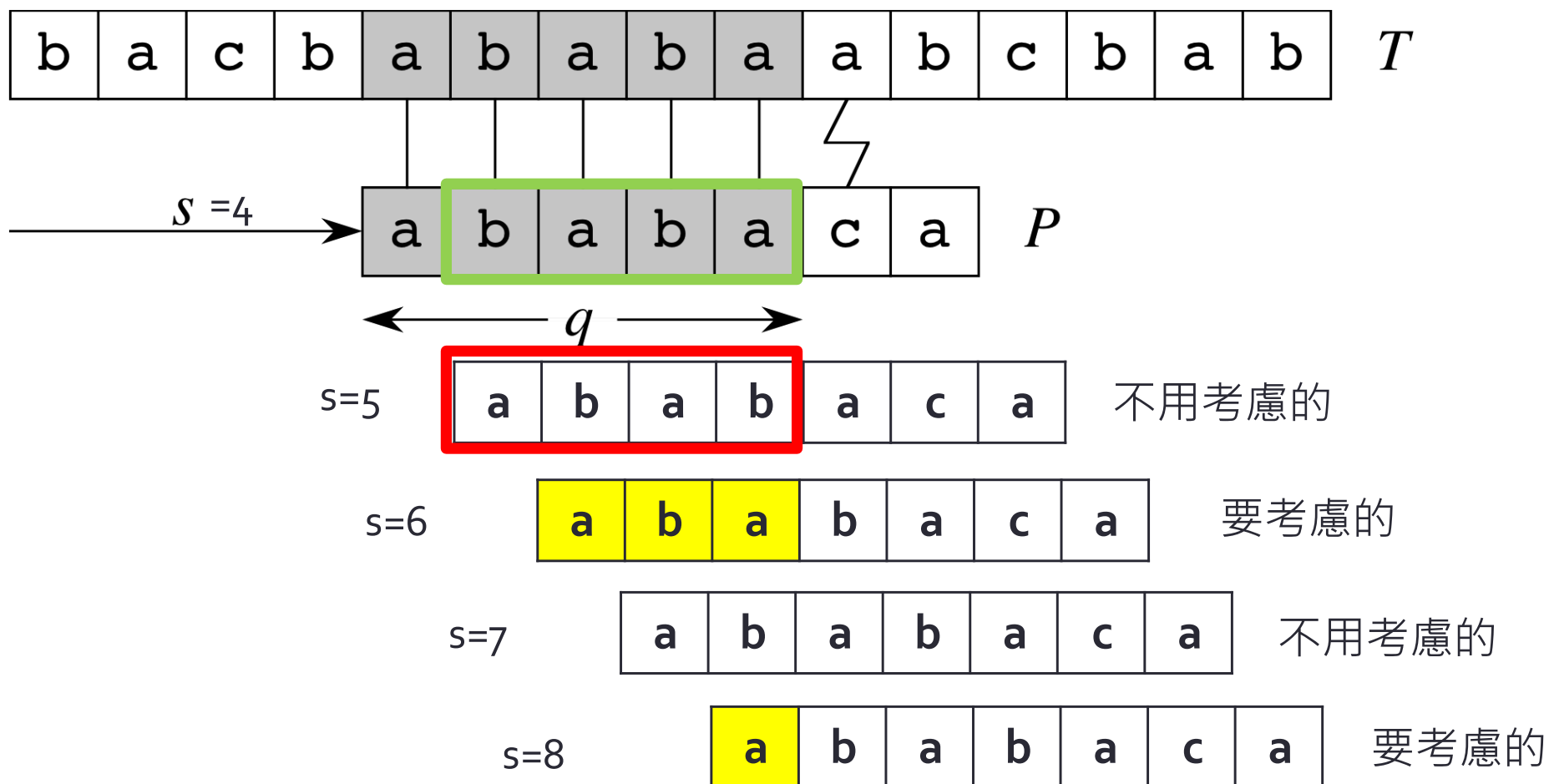
暴力法浪費時間的地方?

- 因為每次for執行比對, 如果錯了, 這回合的資訊完全丟掉.
- 例: P=aaab
- 如果我們發現s=0是valid shift (表示T開頭為aaab),
- 那麼從之前的結果應該可以「知道」shift 1, 2, 3 都可以直接跳過, 不需要一一比對.
- 這個演算法應該要有一些「記憶」

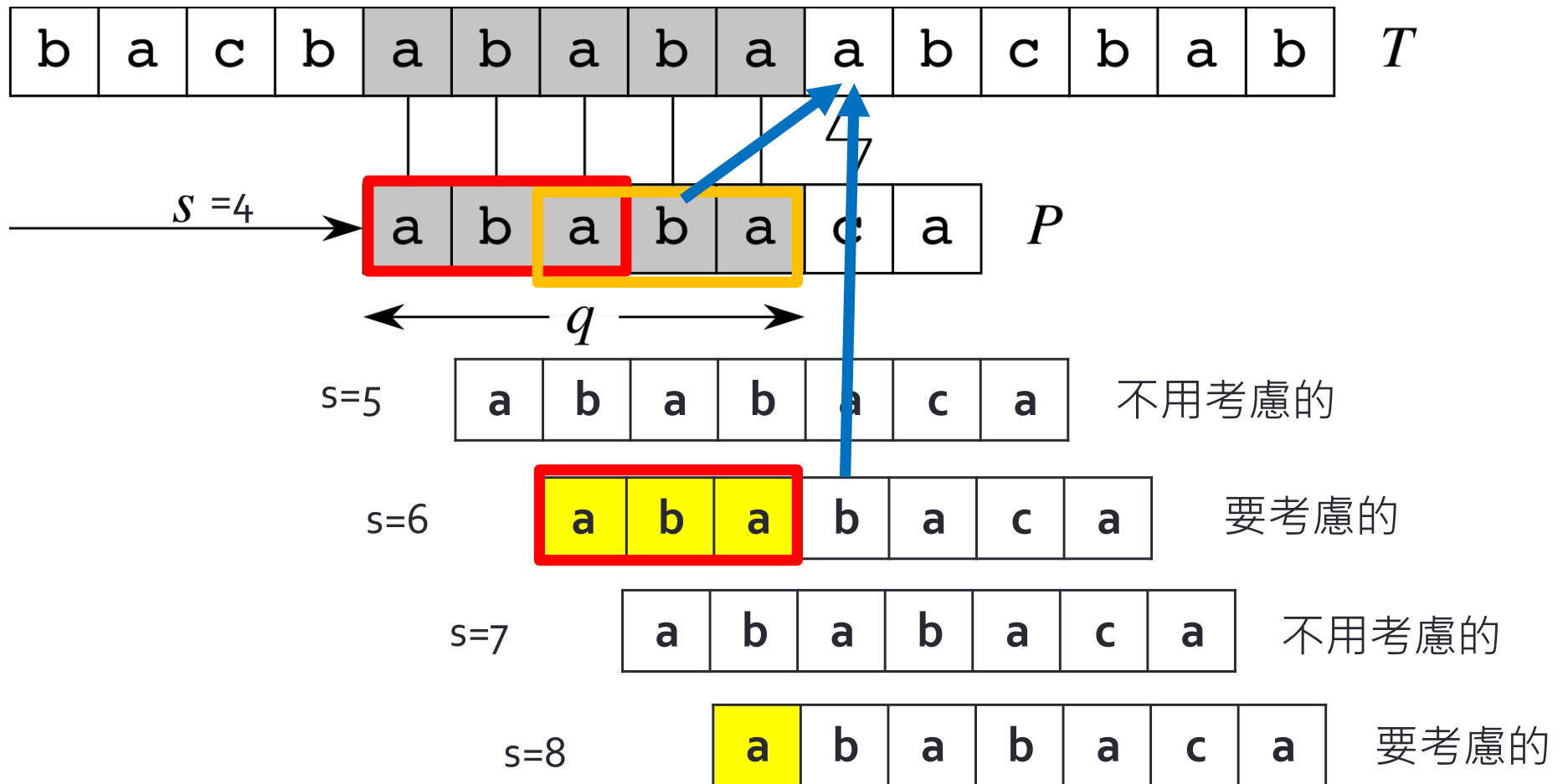


方法二：The Knuth-Morris-Pratt Algo.

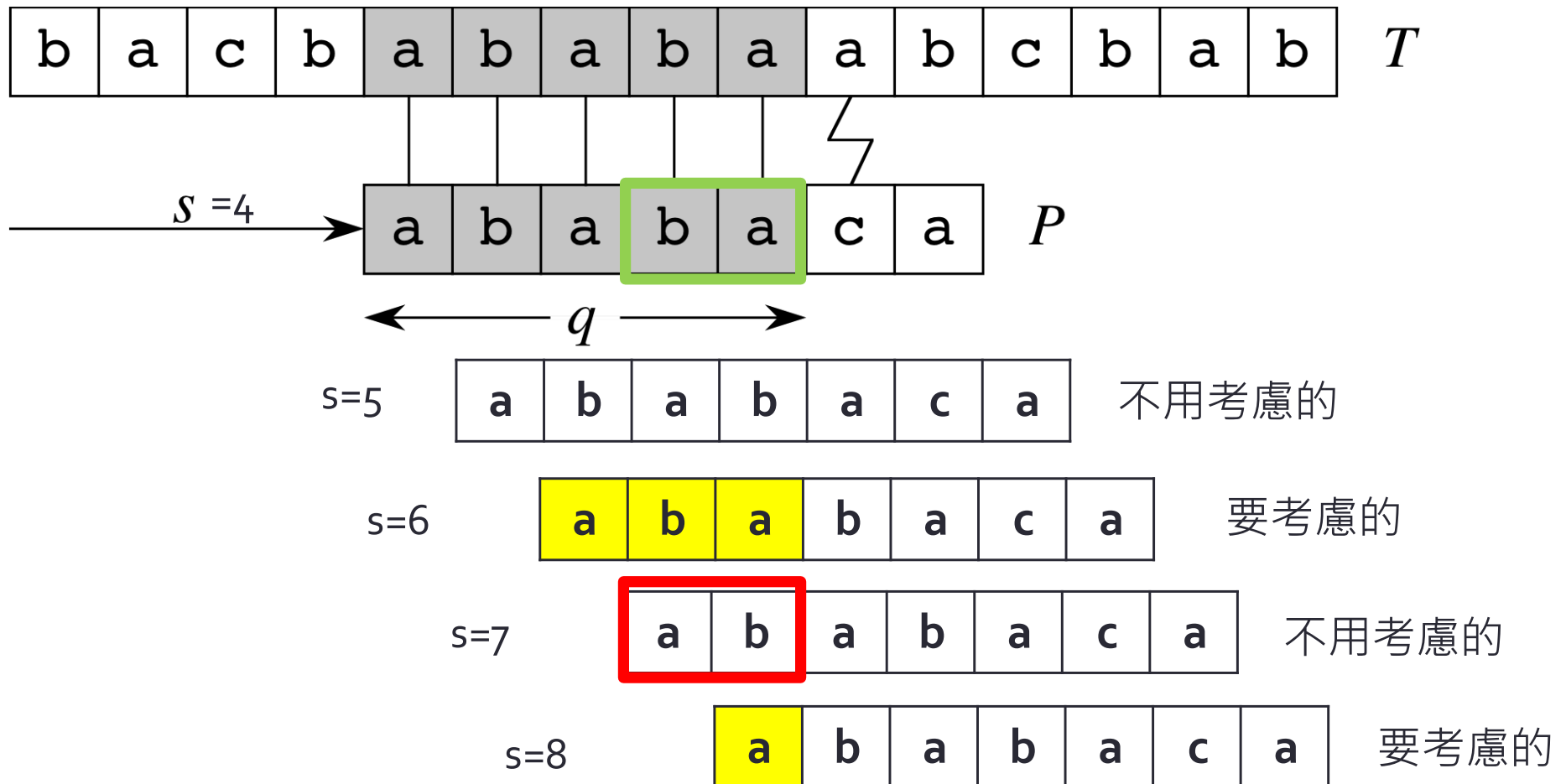




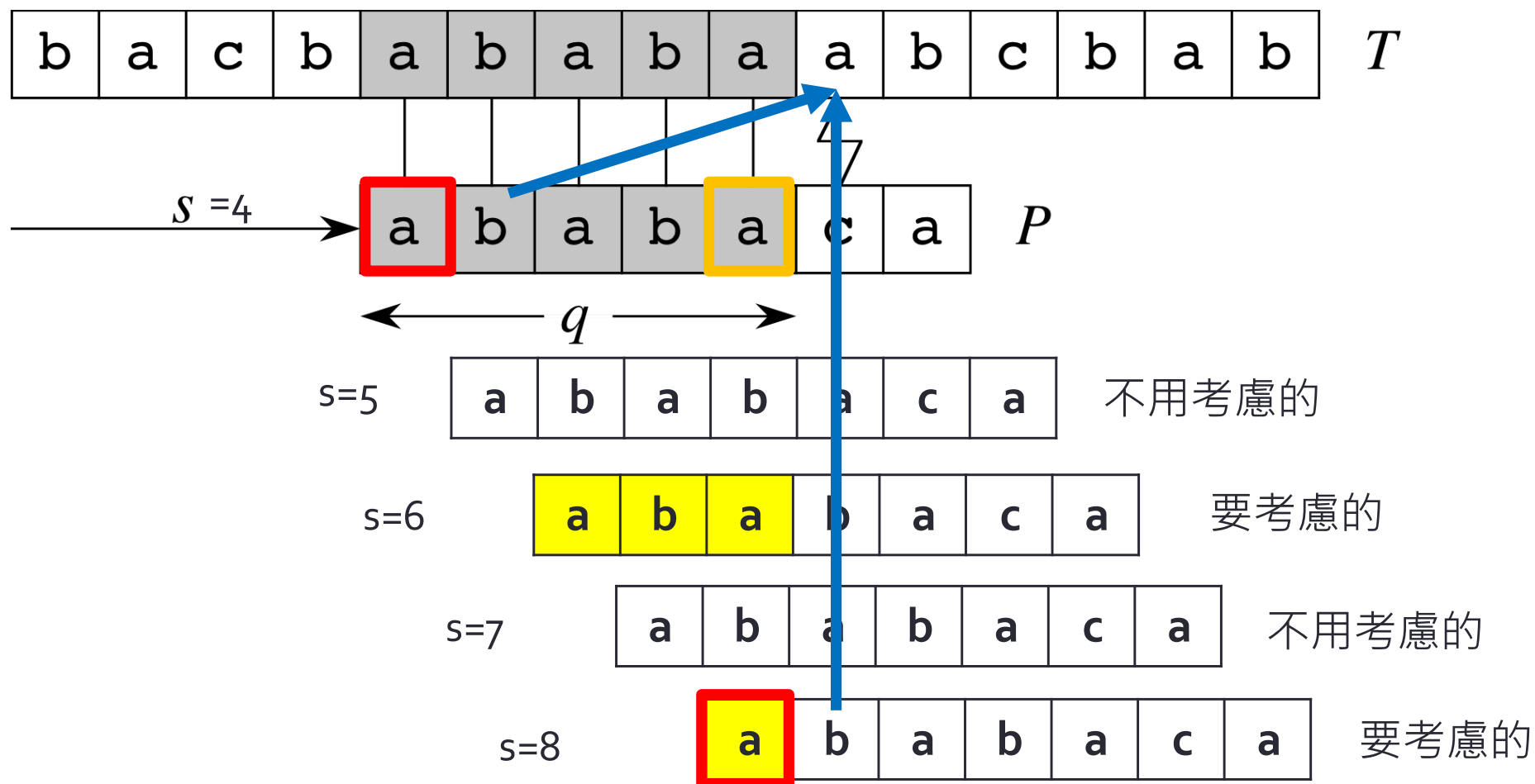
方法二：The Knuth-Morris-Pratt Algo.



方法二：The Knuth-Morris-Pratt Algo.



方法二：The Knuth-Morris-Pratt Algo.



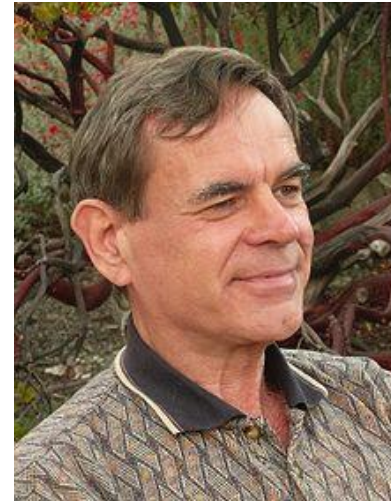
Knuth-Morris-Pratt



Don Knuth



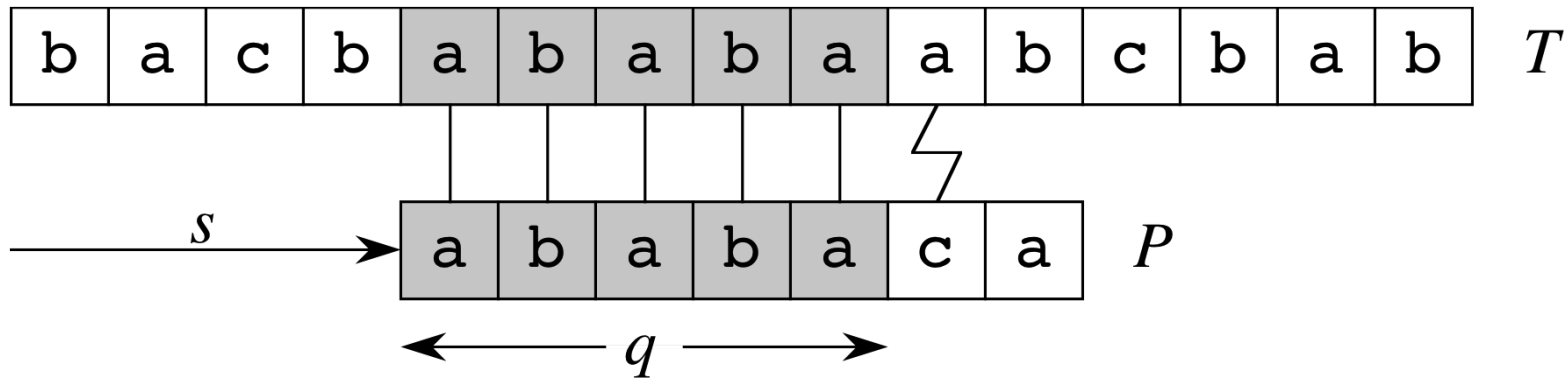
James Morris



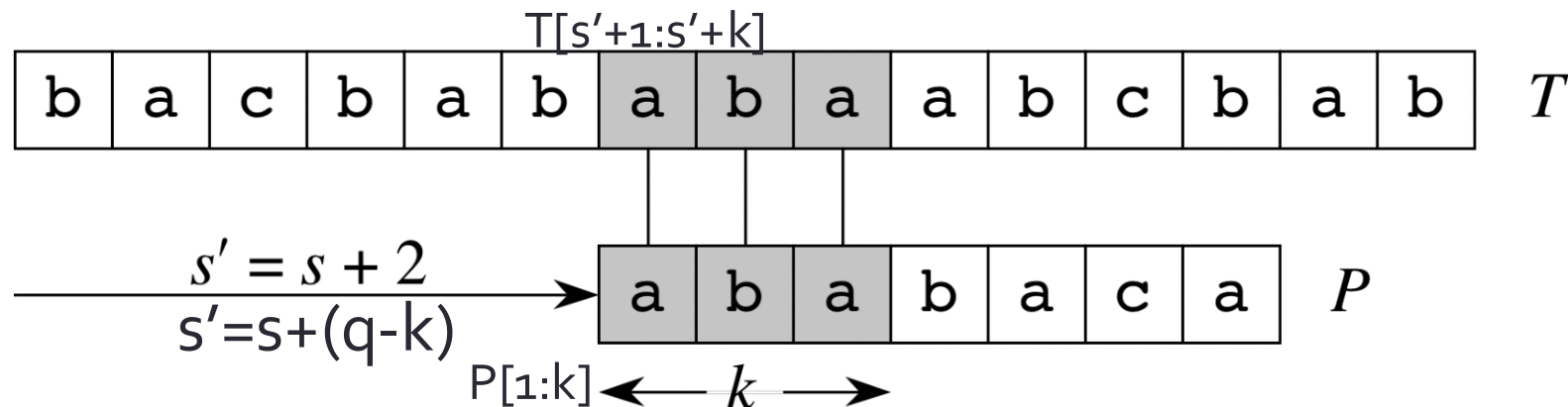
Vaughan Pratt

正式一點的說法

- 假設 $P[1:q]$ 和 $T[s+1:s+q]$ 已經match了



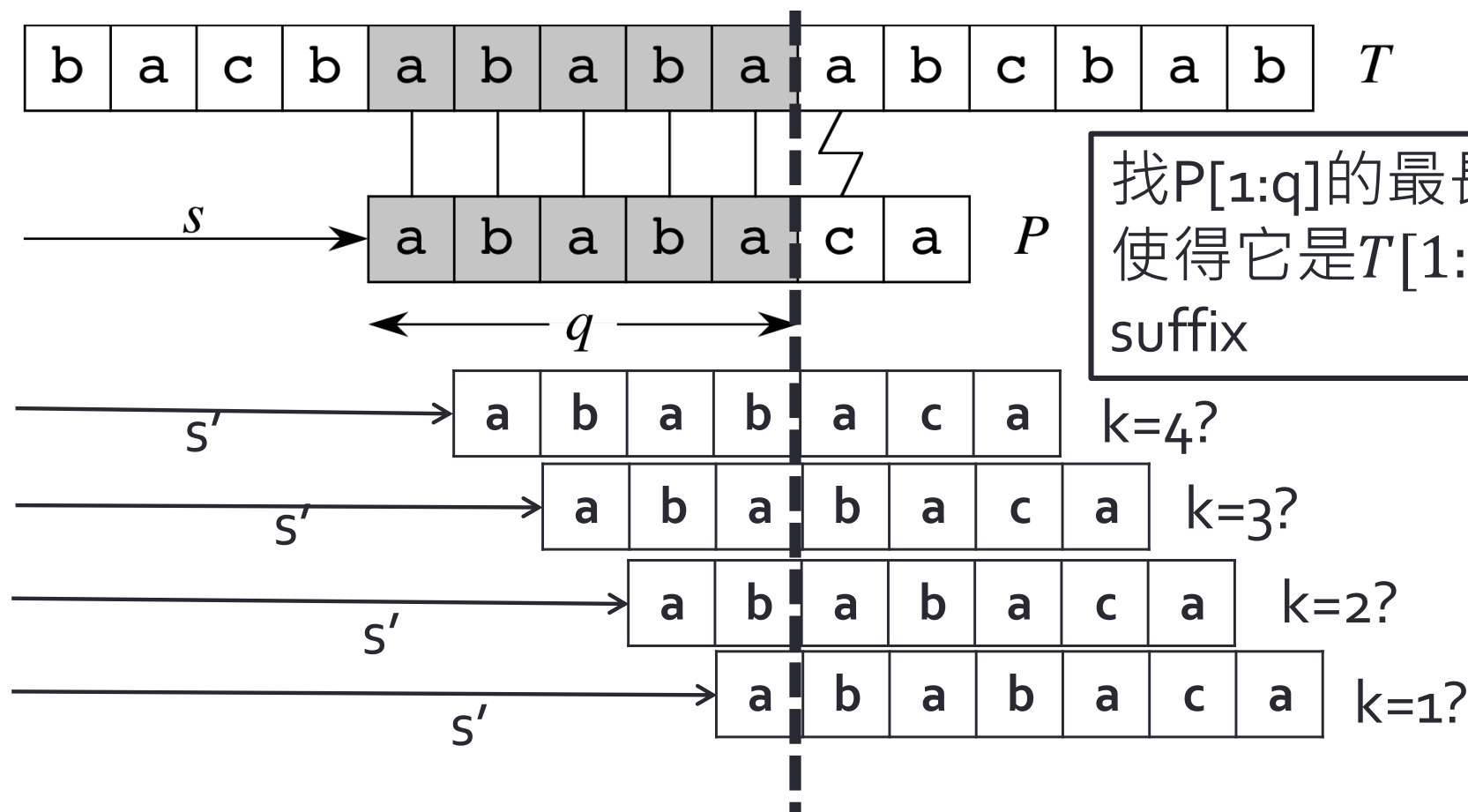
- 要找出最小的shift s' 使得某個 $k < q$ 可以滿足
 $P[1:k] = T[s' + 1:s' + k]$, 且 $s' + k = s + q$



正式一點的說法

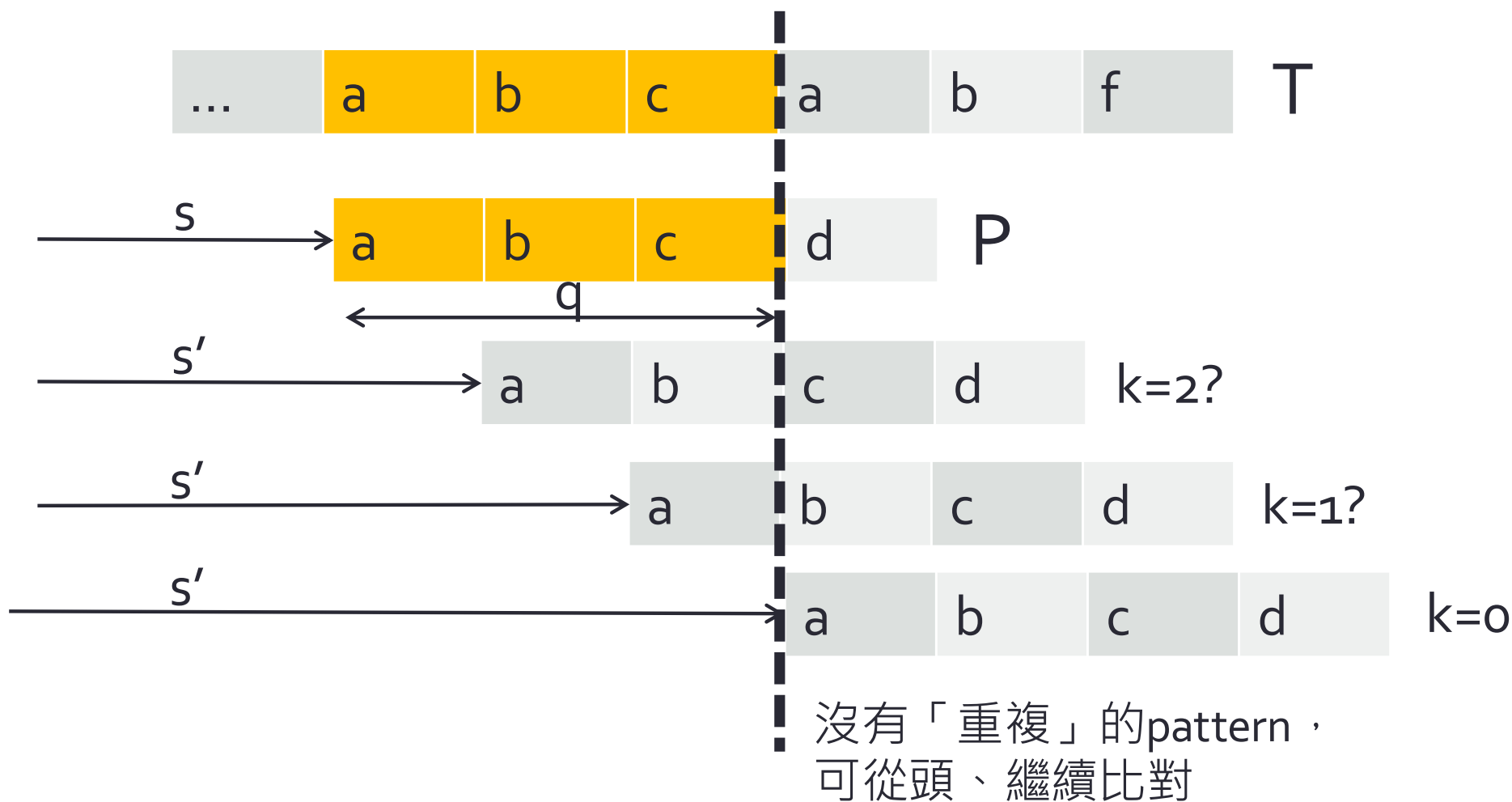
$P[1:q]$ 是 $T[1:s+q]$ 的 suffix

- 假設 $P[1:q]$ 和 $T[s+1:s+q]$ 已經 match 了



- 要找出最小的 shift s' 使得某個 $k < q$ 可以滿足
 $P[1:k] = T[s' + 1:s' + k]$, 且 $s' + k = s + q$

沒有重複的pattern

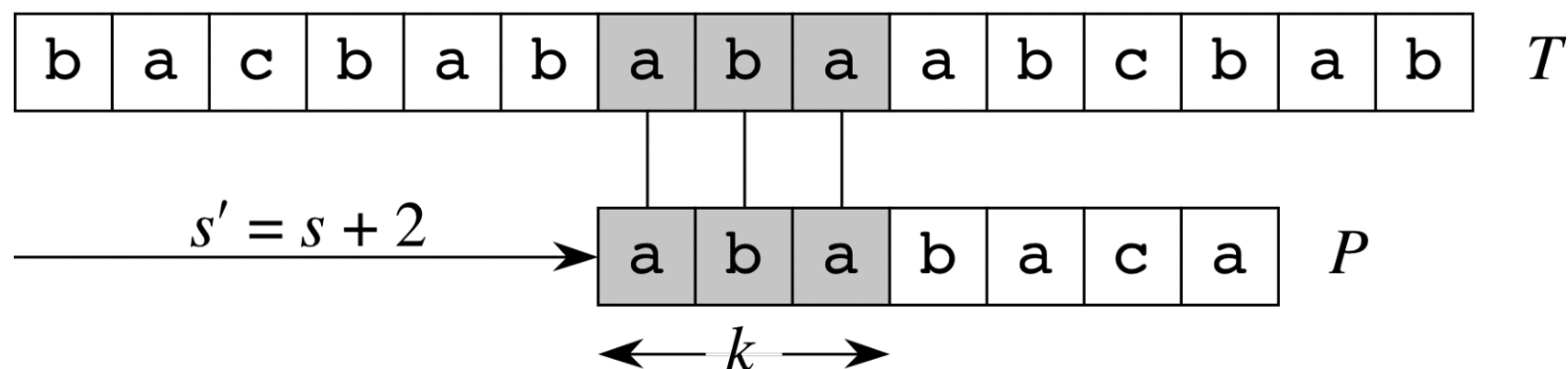
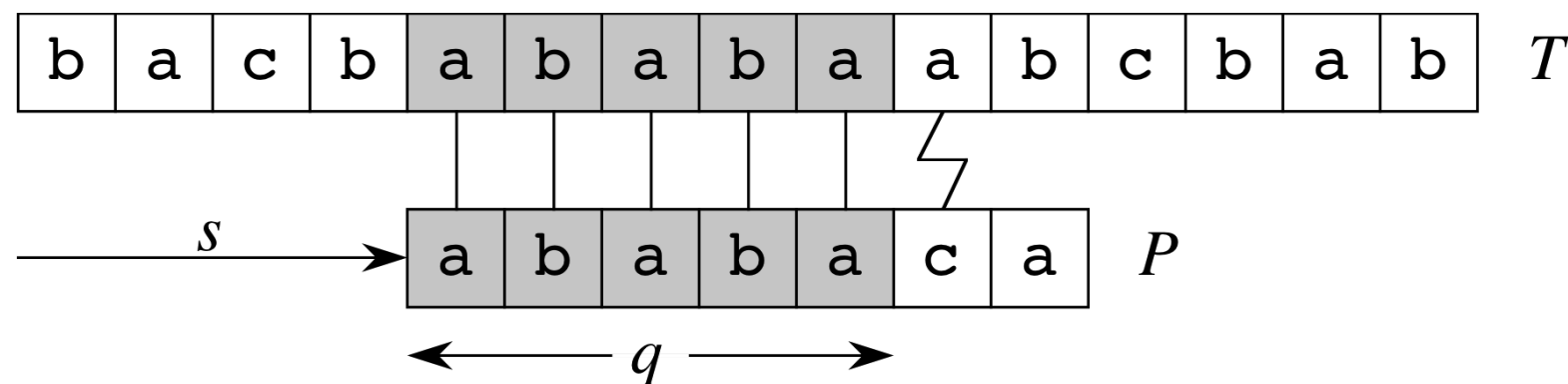


先處理P來取得“重複pattern”的資訊

找 $P[1:q]$ 的最長prefix
使得它是 $T[1:s+q]$ 的suffix



找 $P[1:q]$ 的最長prefix
使得它是 $P[1:q]$ 的suffix



先處理P來取得「重複pattern」的資訊

- 定義Prefix function (failure function) π :
- Input: $\{1, 2, \dots, m\}$
- Output: $\{0, 1, \dots, m-1\}$
- $\pi[q] = \max\{k: k < q \text{ and } P[1:k] \text{ is a suffix of } P[1:q]\}$
- (也就是前面例子中的k值, 可以想成最長的「重複」pattern的長度)

Prefix function example

$$\pi[q] = \max\{k: k < q \text{ and } P_k \text{ is a suffix of } P_q\}$$

i	1	2	3	4	5	6	7
P[i]	A	B	A	B	A	C	A
$\pi[i]$							
$\pi[i]$	0	0	1	2	3	0	1

[illegible]

Pseudo-code: Prefix function

```
1 Compute-Prefix-Function(P)
2 m=P.length
3 let  $\pi[1:m]$  be a new array
4  $\pi[1] = 0$ 
5 k=0
6 for q=2 to m
7     while k>0 and P[k+1] != P[q]
8         k= $\pi[k]$ 
9     if P[k+1]==P[q]
10        k=k+1
11     $\pi[q] = k$ 
12 return  $\pi$ 
```

Pseudo-code: Matcher

```
1 KMP-Matcher (T, P, n, m)
2  $\pi$  = Compute-Prefix-Function
3 q=0
4 k=0
5 for i=1 to n
6     while q>0 and P[q+1]!=T[i]
7         q= $\pi$ [q]
8     if P[q+1]==T[i]
9         q=q+1
10    if q==m
11        print "Pattern occurs with shift" i-m
12        q= $\pi$ [q]
```

例子：Matching

- Ex. 1: $T = \text{BACBABABAABCBAAB}$
- Ex. 2: $T = \text{BABABABABACA}$

i	1	2	3	4	5	6	7
P[i]	A	B	A	B	A	C	A
$\pi[i]$	0	0	1	2	3	0	1

- 請見Cormen p. 978 KMP-Matcher

算Prefix function花多少時間?

Compute-Prefix-Function(P)

$m = P.length$

let $\pi[1:m]$ be a new array

$\pi[1] = 0$

$k = 0$

for $q = 2$ to m 共 $O(m)$ 次

 while $k > 0$ and $P[k+1] \neq P[q]$

$k = \pi[k]$

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

$k = k + 1$

$\pi[q] = k$

return π

麻煩的是這邊：
總共會跳多少次呢?

算Prefix function花多少時間?

Compute-Prefix-Function (P)

$m = P.length$

let $\pi[1..m]$ be a new array

$\pi[1] = 0$

$k = 0$

for $q = 2$ to m

while $k > 0$ and $P[k+1] \neq P[q]$

$k = \pi[k]$

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

$k = k + 1$

$\pi[q] = k$

return π

進入迴圈的時候 $k < q$, 且 q 每次增加, k 有時候不增加
所以 $k < q$ 永遠成立

所以 $\pi[q] = k < q$.

所以每執行一次迴圈
就減少 k 一次
且 k 永遠不是負的

k 只會在這邊增加,
因此最多總共增加 $m-1$ 次 (迴圈執行次數)

最後: 既然有增加才有得減少, while
loop 總共執行的次數不會超過 $O(m)$

Total:
 $O(m)$

Matcher要花多少時間？

```
1 KMP-Matcher (T, P, n, m)
2  $\pi$  = Compute-Prefix-Function
3 q=0
4 k=0
5 for i=1 to n
6     while q>0 and P[q+1]!=T[i]
7         q= $\pi$ [q]
8     if P[q+1]==T[i]
9         q=q+1
10    if q==m
11        print "Pattern occurs with shift" i-m
12        q= $\pi$ [q]
```

KMP執行時間

- 所以總和來看:
- Preprocessing時間 $O(m)$
- 比對時間 $O(n)$



方法三：The Rabin-Karp Algorithm

- 假設 $\Sigma = \{0, 1, \dots, 9\}$
- 那麼每個長度為 k 的字串可以想成是一個 k 位數的十進位數
- 例如字串 "31415" 可以想成是十進位數 31415

checksum

檢查碼

資料



- 把 t_s 設為代表 $T[s+1:s+m]$ 的十進位數
- p 設為代表 P 的十進位數
- 那麼 $t_s = p$ iff $T[s+1:s+m] = P[1:m]$, 也就是 s 是 valid shift

- 怎麼從 P 計算 p 呢?
- $p = P[m] + 10 P[m-1] + 10^2 P[m-2] + \dots + 10^{m-2} P[2] + 10^m P[1]$
- $O(m)$ 的算法
- $p = P[m] + 10(P[m-1] + 10(P[m-2] + \dots + 10(P[2] + 10P[1]) \dots))$

$$P[1] \quad P[2] \quad P[3] \quad \dots$$

$$10^{m-1} \quad 10^{m-2} \quad 10^{m-3}$$

$$P[m-1] \quad P[m]$$

$$10^1 \quad 10^0$$

Q1: 10^k 要算多久? ($0 \leq k \leq m-1$)

$$O()$$

Q2: 整個 P 要算多久?

$$O()$$

1 2 3

另一種算法

$$[10 P[1] + P[2]] \times 10$$

=

$$[(\quad) + P[3]] \times 10$$

$$[\quad] \times 10$$

...

$$\Rightarrow P = 10 \dots 10 (10 (10 P[1] + P[2]) + P[3]) + P[4]) \dots + P[m]$$

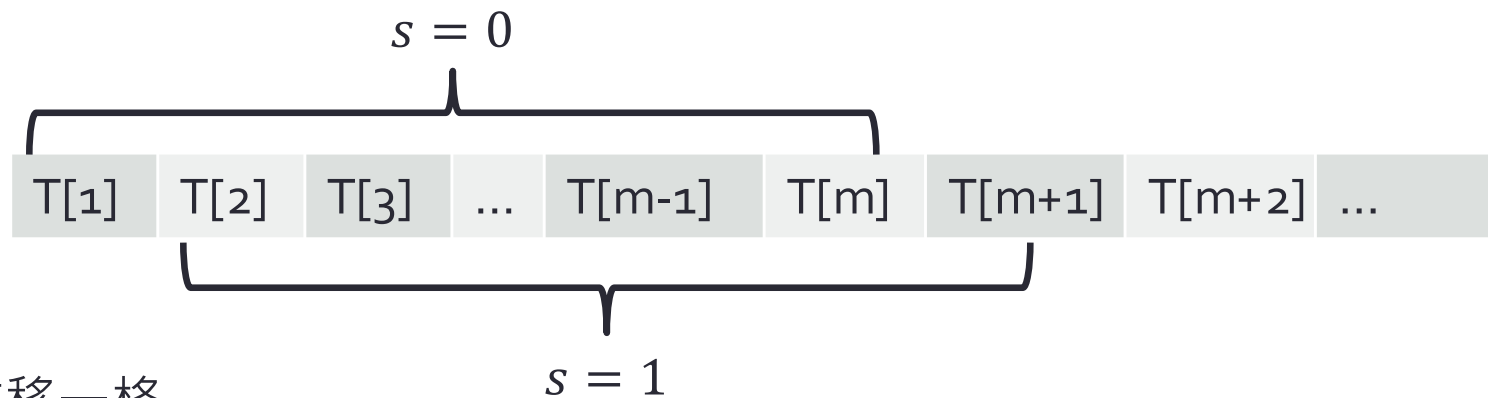
每一項都用一個乘法

一個加法.

m 項 \Rightarrow running time

方法三：The Rabin-Karp Algorithm

- 怎麼從 ~~t_s~~ 計算 ~~t_{s+1}~~ 呢?
- $t_0 = T[m] + 10(T[m-1] + 10(T[m-2] + \dots + 10(T[2] + 10T[1]) \dots))$



整個往右移一格

$$t_{s+1} = 10(t_s - 10^{m-1}T[s+1]) + T[s+m+1]$$

拿掉最左邊那一格

加上最右邊那一格

$$t_s = 10^{m-1} T[s+1] + 10^{m-2} T[s+2] + \dots + 10 T[s+m-1] + T[s+m]$$

$$t_{s+1} = 10^{m-1} T[s+2] + \dots + 10^2 T[s+m-1] + 10 T[s+m] + T[s+m+1]$$

for t_s

$T[s+1]$	$T[s+2]$...	$T[s+m]$	$T[s+m+1]$	String	checksum
					$T[s+1:s+m]$	t_s
10^{m-1}	10^{m-2}	...	10^0			

$t_{s+1} \Rightarrow$

	$T[s+2:s+m+1]$	t_{s+1}
10^{m-1}		
10^1		
10^0		

$$t_s = 10^{m-1} \cdot T[s+1] + \boxed{}$$

pre-calculated in ? time

$$10^{m-1} \boxed{O()}$$

$$t_{s+1} = \boxed{} \cdot 10 + T[s+m+1]$$

$$= 10(t_s - 10^{m-1} \cdot T[s+1]) + T[s+m+1]$$

$$\boxed{O()}$$

方法三：The Rabin-Karp Algorithm

- 那麼用這個方法要多花少時間呢? (簡易分析版)
- $p = P[m] + 10(P[m-1] + 10(P[m-2] + \dots + 10(P[2] + 10P[1]) \dots))$ $O(m)$
- $t_0 = T[m] + 10(T[m-1] + 10(T[m-2] + \dots + 10(T[2] + 10T[1]) \dots))$ $O(m)$
- 然後用 $t_{s+1} = 10(t_s - 10^{m-1}T[s+1]) + T[s+m+1]$, 算 t_1, t_2, \dots, t_{n-m} 的值 (每次都是constant time, 共n-m次)
- 所以總共: $O(m)$ preprocessing時間, $O(n-m)$ 比對時間

方法三：The Rabin-Karp Algorithm

• 之前的兩個問題:

1. Σ 如果是general的character set, 怎麼辦? (不再是 $\{0,1,\dots,9\}$)

如何解決呢?

假設 $|\Sigma| = d$, $\Sigma = \{c_1, c_2, c_3, \dots, c_d\}$. 可以把之前的式子改成

$$\begin{aligned}
 & p \\
 &= \text{idx}(P[m]) \\
 &+ d(\text{idx}(P[m-1]) \\
 &+ d(\text{idx}(P[m-2]) + \dots + d(\text{idx}(P[2]) + d \cdot \text{idx}(P[1])) \dots))
 \end{aligned}$$

(a) 把整個string看成一個d進位的數.

(b) 字元在 Σ 中index當作該字元所代表的值

mini - HW - H

$$\Sigma = \{A, B, C, \dots, Z, a, b, \dots, z\} \quad |\Sigma| = 52$$

index: 0, 1, 2, \dots, 25, 26, 27, \dots, 51

$$P = ABCxyz \quad m=6$$

$$p =$$

$$= 1121919$$

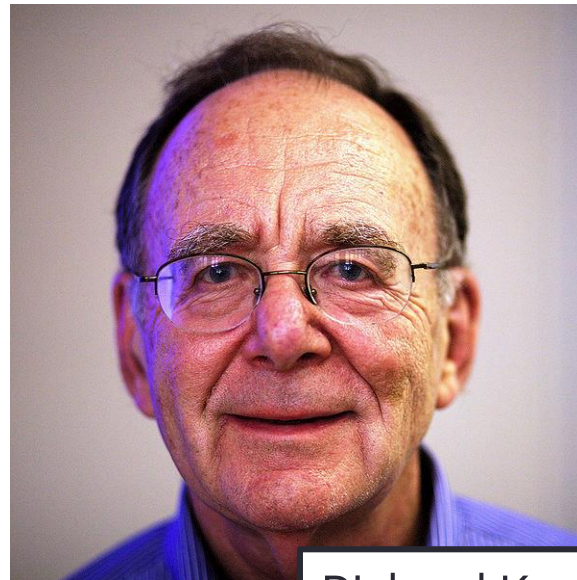
C	idx(C)
A	0
B	1
C	2
x	49
y	50
z	51

方法三：The Rabin-Karp Algorithm

- 2. 當 m 比較大的時候, p 和 t_s 將很難用電腦直接處理 (用long long也存不下)
 - 加一加乘一乘最後總是會overflow
- 如何解決? 利用同餘理論.



Michael Rabin



Richard Karp

同餘理論 (Modular Arithmetic)

- 假設 a, b 都為整數
- $a \equiv b \pmod{n}$: 表示 a 和 b 除以 n 的餘數相等 $0 \leq c < 12$
- 例如:
 - $38 \equiv 14 \pmod{12}$
 - $0 \equiv 12 \pmod{12}$
 - $-13 \equiv -3 \pmod{5}$
- 更棒的性質:
- $a_1 \equiv b_1 \pmod{n}$
- $a_2 \equiv b_2 \pmod{n}$
- 則
 1. $a_1 + a_2 \equiv b_1 + b_2 \pmod{n}$
 2. $a_1 - a_2 \equiv b_1 - b_2 \pmod{n}$
 3. $a_1 a_2 \equiv b_1 b_2 \pmod{n}$

① what if in $a \pmod n$ $a < 0$?

the least positive number
that should be subtracted
from the dividend to make it
divisible.

$$-13 \pmod 5 \rightarrow -13 - 2 = -15 \text{ divisible by } 5$$

$$-3 \pmod 5 \rightarrow -3 - 2 = -5 \text{ divisible by } 5$$

② $a_1 \equiv b_1 \pmod n$

$$a_1 + a_2 \equiv b_1 + b_2 \pmod n$$

$$a_2 \equiv b_2 \pmod n$$

沒 mod 的 a_1 a_2 mod 的 b_1 b_2

$$17 \pmod 5 = 2$$

$$(17+8) \pmod 5 = 0$$

$$8 \pmod 5 = 3$$

$$(2+3) \pmod 5 = 0$$

實際計算時，先 mod 比較好算。

同餘理論 (Modular Arithmetic)

- $a_1 \equiv b_1 \pmod{n}$
- $a_2 \equiv b_2 \pmod{n}$
- 則
- $a_1 + a_2 \equiv b_1 + b_2 \pmod{n}$
- 證明:
- $a_1 \equiv b_1 \pmod{n}$ 表示
 - $a_1 = c_{a_1}n + r_1$
 - $b_1 = c_{b_1}n + r_1$
- $a_2 \equiv b_2 \pmod{n}$ 表示
 - $a_2 = c_{a_2}n + r_2$
 - $b_2 = c_{b_2}n + r_2$
- 所以
 - $a_1 + a_2 = (c_{a_1} + c_{a_2})n + (r_1 + r_2)$
 - $b_1 + b_2 = (c_{b_1} + c_{b_2})n + (r_1 + r_2)$
- 兩者餘數相同!
- $a_1 + a_2 \equiv b_1 + b_2 \pmod{n}$ 得證.

同餘理論 (Modular Arithmetic)

- $a_1 \equiv b_1 \pmod{n}$
- $a_2 \equiv b_2 \pmod{n}$
- 則
- $a_1 a_2 \equiv b_1 b_2 \pmod{n}$
- 證明:
- $a_1 \equiv b_1 \pmod{n}$ 表示
 - $a_1 = c_{a_1} n + r_1$
 - $b_1 = c_{b_1} n + r_1$
- $a_2 \equiv b_2 \pmod{n}$ 表示
 - $a_2 = c_{a_2} n + r_2$
 - $b_2 = c_{b_2} n + r_2$
- 所以
 - $a_1 a_2$

$$= (c_{a_1} n + r_1)(c_{a_2} n + r_2)$$

$$= (c_{a_1} c_{a_2} n + c_{a_1} r_2 + c_{a_2} r_1) n + r_1 r_2$$
 - $b_1 b_2 = (c_{b_1} c_{b_2} n + c_{b_1} r_2 + c_{b_2} r_1) n + r_1 r_2$
- 兩者餘數相同! 得證!

The Rabin-Karp Algorithm 修正版

- 取 q 使得 dq 可以用一個電腦word (32-bit or 64-bit)來表示
- 既然mod後再加, 減, 乘也會保持原本的關係, 我們可以把這些operation都變成mod版本的
- $p = P[m] + d(P[m-1] + d(P[m-2] + \dots + d(P[2] + dP[1]) \dots)) \pmod{q}$
- $t_0 = T[m] + d(T[m-1] + d(T[m-2] + \dots + d(T[2] + dT[1]) \dots)) \pmod{q}$
- $t_{s+1} = d(t_s - d^{m-1}T[s+1]) + T[s+m+1] \pmod{q}$

The Rabin-Karp Algorithm 修正版

- 新的mod版algorithm會造成一個問題:

- 雖然 $t_s = p \Rightarrow t_s \equiv p \pmod{q}$
- 但 $t_s = p \not\Leftrightarrow t_s \equiv p \pmod{q}$
- 例如 $38 \equiv 14 \pmod{12}$, $38 \neq 14$
- 但是如果 $t_s \not\equiv p \pmod{q} \Rightarrow t_s \neq p$

data1 \searrow Same
data2 \rightarrow checksum

- 所以演算法變成這樣:

1. 如果 $t_s \not\equiv p \pmod{q}$, 那麼現在這個s為invalid shift
2. 如果 $t_s \equiv p \pmod{q}$, 那麼必須額外檢查
(直接比對範圍內的字串 \rightarrow 花很多時間)

- 當 $t_s \equiv p \pmod{q}$, 但是 $t_s \neq p$ 時, 稱為spurious hit
- 當q夠大的時候, 希望spurious hit會相當少

為什麼q最好和d互質? $d = |\Sigma|$

希望checksum越分散越好

q最好是質數, 因為d不確定

比如說 $d=10$ $q=500$

$$T[S+1] \times 10^{m-1} + \dots + T[S+m-4] \times 10^4 +$$

$$T[S+m-3] \times 10^3 + T[S+m-2] \times 10^2 + T[S+m-1] \times 10$$

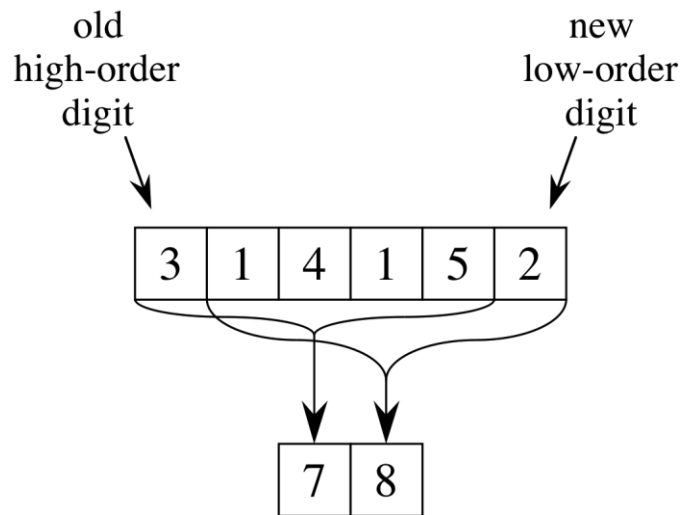
$$+ T[S+m]$$

$$(\text{mod } q)$$

不會被 checksum 考慮.

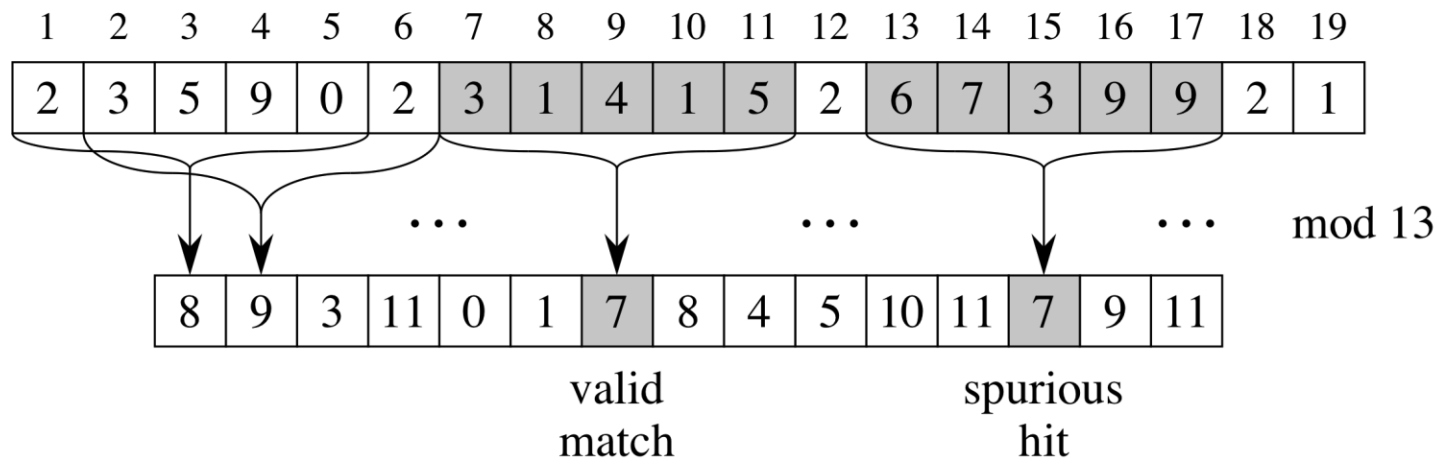
$$(\text{mod } q)$$

例子: Rabin-Karp Algorithm



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

Diagram illustrating the rolling hash calculation for a sequence of digits. The input sequence is 2 3 5 9 0 2 3 1 4 1 5 2 6 7 3 9 9 2 1. The first digit (2) is labeled "old high-order digit", the last digit (1) is labeled "new low-order digit", and the shift operation is indicated by an arrow.



Pseudo Code: Rabin-Karp

Rabin-Karp-Matcher (T, P, d, q)

n=T.length

m=P.length

$h = d^{m-1} \bmod q$

p=0

t=0

for i=1 to m

$p = (dp + P[i]) \bmod q$

$t = (dt + T[i]) \bmod q$

for s=0 to n-m

迴圈跑n-m+1次

 if p==t

 if $P[1..m] == T[s+1..s+m]$

Hit的時候比對: $O(m)$

 print "Pattern occurs with shift" s

 if s<n-m

$t = (d(t - T[s+1]h) + T[s+m+1]) \bmod q$

T: string to be searched
P: pattern to be matched
d: size of the character set
q: max number

Pre-processing: $O(m)$

Worst-case Running Time

Rabin-Karp-Matcher (T, P, d, q)

$n = T.length$

$m = P.length$

$h = d^{m-1} \bmod q$

$p = 0$

$t = 0$

for $i = 1$ to m

$p = (dp + P[i]) \bmod q$

$t = (dt + T[i]) \bmod q$

for $s = 0$ to $n - m$

迴圈跑 $n - m + 1$ 次

if $p == t$

if $P[1..m] == T[s+1..s+m]$

Hit的時候比對: $O(m)$

print "Pattern occurs with shift" s

if $s < n - m$

$t = (d(t - T[s+1]h) + T[s+m+1]) \bmod q$

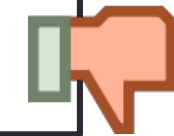
Worst case的時候:

$T = a^n$ (n個a)

$P = a^m$ (m個a)

比對的時間為

$O(m(n-m+1))$



$O(nm)$

Pre-processing: $O(m)$

Average Running Time

```
Rabin-Karp-Matcher (T, P, d, q)
n=T.length
m=P.length
h= $d^{m-1} \bmod q$ 
p=0
```

Pre-processing: $O(m)$

```
for i=1 to m
```

```
    p=(dp+P[i]) mod q
```

迴圈跑 $n-m+1$ 次

```
    (dt+T[i]) mod q
```

```
for s=0 to n-m
```

```
    if p==t
```

Hit的時候比對: $O(m)$

```
    if P[1..m]==T[s+1..s+m]
```

```
        print "Pattern occurs with shift" s
```

```
    if s<n-m
```

```
        t=(d(t-T[s+1]h)+T[s+m+1]) mod q
```

平常的時候, valid shift 很少

(假設有 c 個)

不會每次都有 modulo 的 hit.

假設字串各種排列組合出現的機率相等
則 spurious hit 的機率可當成 $1/q$.

則比對花的時間:

Spurious hit 共花 $O((n-m+1)/q) = O(n/q)$ 次

總共比對花的時間為

$O((n-m+1) + (m(c+n/q)))$

If $c=O(1)$ and $q \geq m$, $\rightarrow O(n+m) = O(n)$



Assumption:

c 個 valid shifts $C = O(1)$

$\frac{1}{q}$ spurious hit 機率 $q \geq m$
(不含 valid shift 機率)

valid shift:

額外花 $O(cm)$ 比對

Spurious hit:

共有 $n-m+1$ 可能的 shift values

Spurious hit = 次數 $\frac{n-m+1}{q}$

額外花 $O(\frac{n-m+1}{q}m)$

比對 total:

$$O(m(c + \frac{n-m+1}{q}))$$

$$= O(m + n - m + 1) = O(m + n)$$

Related Reading

- Textbook (Cormen)
ch. 32, 32.1, 32.2, 32.4 (正確性的證明略為複雜)