



哈尔滨工业大学

海量数据计算研究中心

Massive Data Computing Lab @ HIT

# 算法设计与分析—入门篇

## 第七讲 字符串搜索

哈尔滨工业大学

王宏志

[wangzh@hit.edu.cn](mailto:wangzh@hit.edu.cn)

<http://homepage.hit.edu.cn/pages/wang/>

# 提纲

**7.1 字符串搜索概述**

**7.2 Rabin-Karp算法**

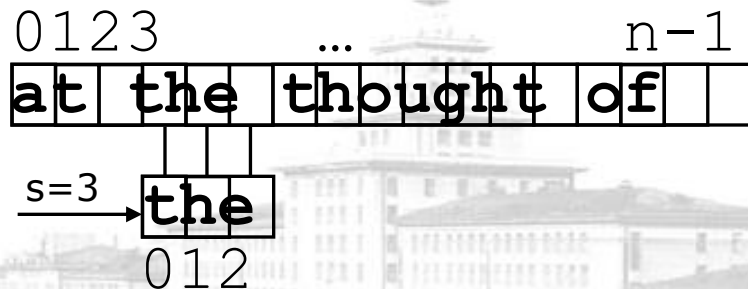
**7.3 KMP算法**

**7.4 BMH算法**



# 字符串匹配问题

- 输入:
  - 文本  $T = \text{"at the thought of"}$ 
    - $n = \text{length}(T) = 17$
  - 模式  $P = \text{"the"}$ 
    - $m = \text{length}(P) = 3$
- 输出:
  - 移动到  $s$  – 最小的整数 ( $0 \leq s \leq n - m$ ) 满足  $T[s .. s+m-1] = P[0 .. m-1]$ . 返回  $-1$ , 如果不存在这样的  $s$



# 简单匹配算法

- 想法：暴力搜索
  - 检查从0 到  $n - m$ 的所有值

**Naive-Search** ( $T, P$ )

```
01 for  $s \leftarrow 0$  to  $n - m$ 
02      $j \leftarrow 0$ 
03     // check if  $T[s..s+m-1] = P[0..m-1]$ 
04     while  $T[s+j] = P[j]$  do
05          $j \leftarrow j + 1$ 
06         if  $j = m$  return  $s$ 
07 return  $-1$ 
```

- 令  $T = \text{"at the thought of"}$ ,  $P = \text{"though"}$ 
  - 需要多少次比较?

# 简单匹配算法的分析

- 最坏情况:
  - 外层循环:  $n - m$
  - 内层循环:  $m$
  - 总计  $(n-m)m = O(nm)$
  - 何种输入产生最坏情况?
- 最好情况:  $n-m$ 
  - 何时?
- 完全随机的文本和模式:
  - $O(n-m)$

# 提纲

7.1 字符串搜索概述

7.2 Rabin-Karp算法

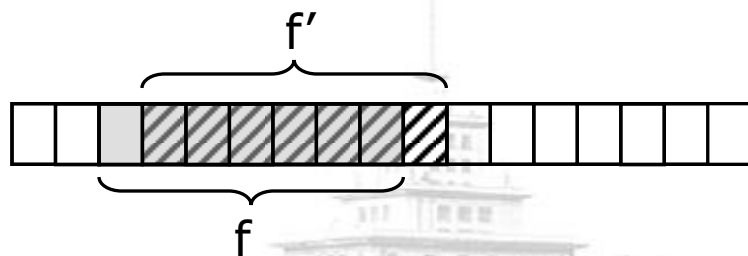
7.3 KMP算法

7.4 BMH算法



# 指纹想法

- 假设:
  - 我们可以在 $O(m)$ 时间计算一个 $P$ 的指纹 $f(P)$ .
  - 如果  $f(P) \neq f(T[s .. s+m-1])$ , 那么  $P \neq T[s .. s+m-1]$
  - 我们可以在 $O(1)$ 时间比较指纹
  - 我们可以在 $O(1)$ 的时间从 $f(T[s .. s+m-1])$ 计算 $f' = f(T[s+1 .. s+m])$

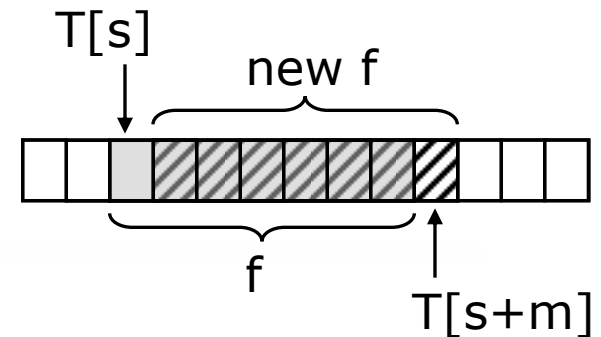


# 基于指纹的算法

- 令字母表位  $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$
- 令指纹为一个十进制数, 即,  $f(\text{"1045"}) = 1 \cdot 10^3 + 0 \cdot 10^2 + 4 \cdot 10^1 + 5 = 1045$

**Fingerprint-Search** ( $T, P$ )

```
01 fp  $\leftarrow$  compute f(P)
02 f  $\leftarrow$  compute f(T[0..m-1])
03 for s  $\leftarrow$  0 to n - m do
04     if fp = f return s
05     f  $\leftarrow$  (f - T[s] * 10m-1) * 10 + T[s+m]
06 return -1
```



■ 运行时间是  $2O(m) + O(n-m) = O(n)$ !

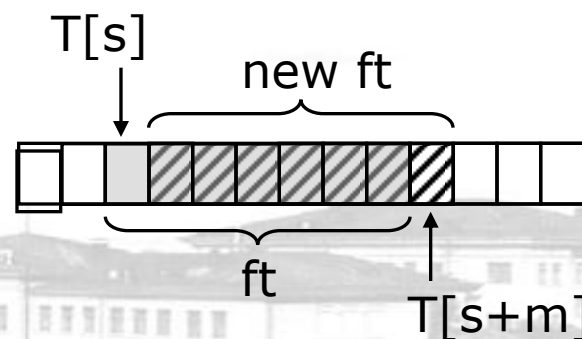


# 使用Hash函数

- 问题: 我们不能假设我们可以对 $m$ 位数在 $O(1)$ 时间内进行算术运算
- 解决方案: 使用hash函数  $h = f \bmod q$ 
  - 例如, 如果  $q = 7$ ,  $h(\text{"52"}) = 52 \bmod 7 = 3$
  - $h(S_1) \neq h(S_2) \Rightarrow S_1 \neq S_2$
  - 但  $h(S_1) = h(S_2)$  并不意味着  $S_1 = S_2$ !
    - 例如, 如果  $q = 7$ ,  $h(\text{"73"}) = 3$ , 但  $\text{"73"} \neq \text{"52"}$
- 但 “ $\bmod q$ ” 算术运算:
  - $(a+b) \bmod q = (a \bmod q + b \bmod q) \bmod q$
  - $(a*b) \bmod q = (a \bmod q) * (b \bmod q) \bmod q$

# 预处理与步骤

- 预处理:
  - $fp = P[m-1] + 10*(P[m-2] + 10*(P[m-3] + \dots + 10*(P[1] + 10*P[0])\dots)) \bmod q$
  - 同样地可以从  $T[0..m-1]$  计算  $ft$
  - 例如:  $P = \text{"2531"} , q = 7, fp$  是多少?
- 步骤:
  - $ft = (ft - T[s]*10^{m-1} \bmod q)*10 + T[s+m]) \bmod q$
  - $10^{m-1} \bmod q$  在预处理中计算一次
  - 例: Let  $T[\dots] = \text{"5319"} , q = 7$ , 对应的  $ft$  是多少?



# Rabin-Karp算法

## Rabin-Karp-Search (T, P)

```
01 q ← a prime larger than m
02 c ←  $10^{m-1} \bmod q$  // run a loop multiplying by 10 mod q
03 fp ← 0; ft ← 0
04 for i ← 0 to m-1 // preprocessing
05     fp ← (10*fp + P[i]) mod q
06     ft ← (10*ft + T[i]) mod q
07 for s ← 0 to n - m // matching
08     if fp = ft then // run a loop to compare strings
09         if P[0..m-1] = T[s..s+m-1] return s
10     ft ← ((ft - T[s]*c)*10 + T[s+m]) mod q
11 return -1
```

- 如果T = "2531978", P = "1978"需要比较字符多少次?

# 分析

- 如果  $q$  是素数, hash函数将会使 $m$ 位字符串在 $q$ 个值中均匀分配
  - 因此, 仅有 $s$ 个轮换中的每第 $q$ 次才需要匹配指纹 (匹配需要比较 $O(m)$ 次)
- 期望运行时间 (如果  $q > m$ ):
  - 预处理:  $O(m)$
  - 外循环:  $O(n-m)$
  - 所有内循环:  $\frac{n-m}{q} m = O(n-m)$
  - 总时间:  $O(n-m)$
- 最坏运行时间:  $O(nm)$

# 应用中的Rabin-Karp算法

- 如果字母表有 $d$ 个字母，将字母翻译为 $d$ 进制数字，即用 $d$ 代替算法中的10
- 选择素数 $q > m$  可以利用随机算法在 $O(m)$ 时间内完成, 或者 $q$  是一个固定大素数且一个计算机字可以容纳 $10 * q$ .
- Rabin-Karp比较简单，可以容易地拓展到2维模式匹配.

# 提纲

**7.1 字符串搜索概述**

**7.2 Rabin-Karp算法**

**7.3 KMP算法**

**7.4 BMH算法**



# n次比较的匹配

- 目标:文本中的每个字符仅匹配一次!
- 简单算法的问题:
  - 没有利用已有部分匹配中的知识
  - 例:

- $T = \text{"Tweedledee and Tweedledum"}$   
     $\text{"Tweedledum"}$

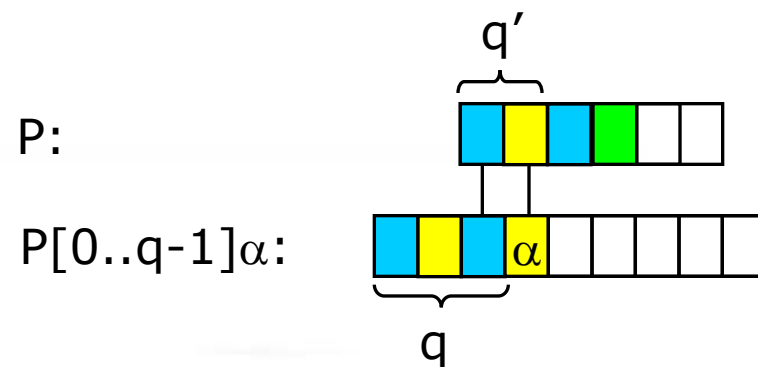
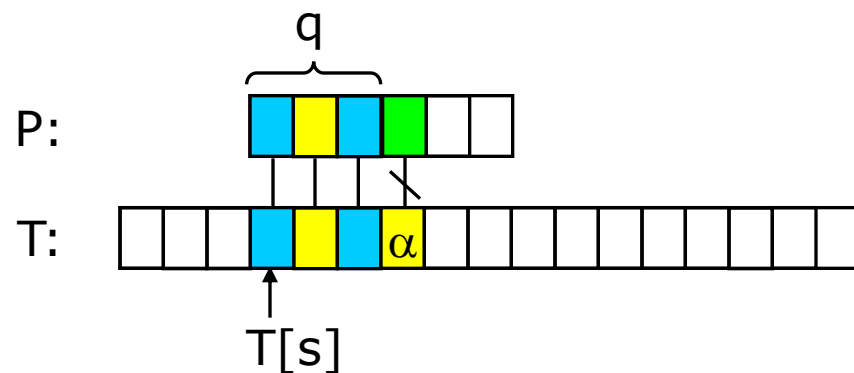
$P =$

- $T = \text{"pappar"}$

$P = \text{"pappappappar"}$

# 一般情况

- 算法的状态:
  - 检查移动  $s$ ,
  - 匹配了  $P$  中的  $q$  个字符
  - 在  $T$  中看到了一个未匹配的字符  $\alpha$ .
- 需要寻找:
  - 最大前缀 “ $P$ -” 并且是  $P[1..q-1]\alpha$  的后缀:



- $q' = \max\{k \leq q \mid P[1..k] = P[q-k+1..q]\alpha\}$

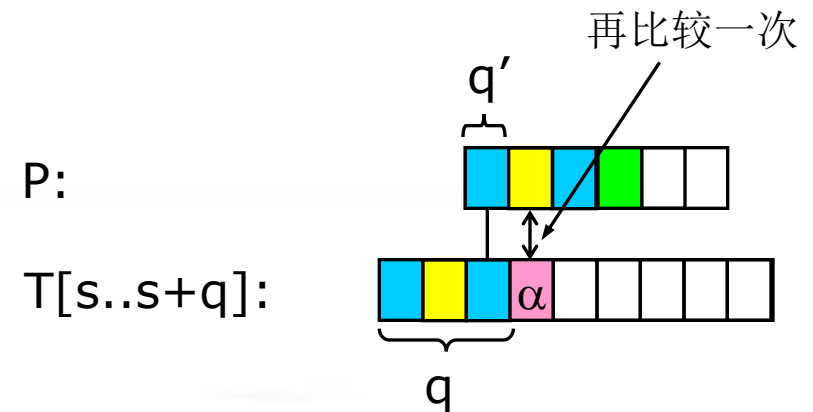
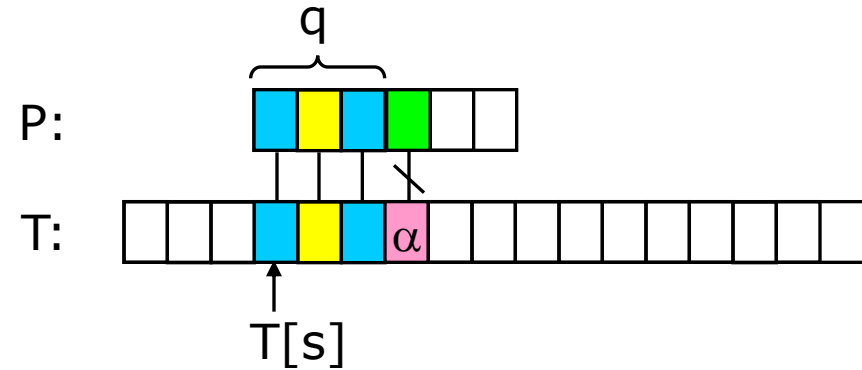


# 自动机搜索

- 算法:
  - 预处理:
    - 对于每个  $q$  ( $1 \leq q \leq m-1$ ) 和每个  $\alpha \in \Sigma$  预先计算一个  $q$  的新值, 记为  $\sigma(q, \alpha)$
    - 填一个大小为  $m/|\Sigma|$  的表
  - 扫描文本
    - 当不匹配发现时 ( $P[q] \neq T[s+q]$ ):
    - 置  $s = s + q - \sigma(q, \alpha) + 1$
    - 分析:
  - ☺ 匹配阶段  $O(n)$
  - ☹ 内存过多:  $O(m/|\Sigma|)$ , 过多的预处理  $O(m/|\Sigma|)$ .

# 前缀函数

- Idea: 忘记未匹配的字符 ( $\alpha$ )!
- 算法的状态:
  - 检查变换  $s$ ,
  - 匹配了  $P$  中的  $q$  个字符
  - 发现了  $T$  中不匹配的字符  $\alpha$ .
- 需要发现:
  - 最大前缀 “ $P$ -” 并且是  $P[1..q-1]$  的后缀:



- $q' = \pi[q] = \max\{k < q \mid P[1..k] = P[q-k+1..q]\}$

# 前缀表

- 我们可以预先计算大小为 $m$ 的前缀表来存储 $\pi[q]$ 的值 ( $0 \leq q < m$ )

$P$		<b>p</b>	<b>a</b>	<b>p</b>	<b>p</b>	<b>a</b>	<b>r</b>
$q$	0	1	2	3	4	5	6
$\pi[q]$	0	0	0	1	1	2	0

- 计算 $P = \text{"dadadu"}$  的前缀表

# Knuth-Morris-Pratt 算法

**KMP-Search**(T, P)

```
01  $\pi \leftarrow \text{Compute-Prefix-Table}(P)$ 
02  $q \leftarrow 0$  // number of characters matched
03 for  $i \leftarrow 0$  to  $n-1$  // scan the text from left to right
04     while  $q > 0$  and  $P[q] \neq T[i]$  do
05          $q \leftarrow \pi[q]$ 
06     if  $P[q] = T[i]$  then  $q \leftarrow q + 1$ 
07     if  $q = m$  then return  $i - m + 1$ 
08 return  $-1$ 
```

- **Compute-Prefix-Table**是P上执行KMP算法的本质.

# KMP的分析

- 最坏运行时间:  $O(n+m)$ 
  - 主算法:  $O(n)$
  - **Compute-Prefix-Table**:  $O(m)$
- 空间:  $O(m)$

# 提纲

**7.1 字符串搜索概述**

**7.2 Rabin-Karp算法**

**7.3 KMP算法**

**7.4 BMH算法**



# 逆简单算法

- 如果从 $P$ 的后面开始搜索？
  - Boyer and Moore

**Reverse-Naive-Search** ( $T, P$ )

```
01 for  $s \leftarrow 0$  to  $n - m$ 
02      $j \leftarrow m - 1$     // start from the end
03     // check if  $T[s..s+m-1] = P[0..m-1]$ 
04     while  $T[s+j] = P[j]$  do
05          $j \leftarrow j - 1$ 
06         if  $j < 0$  return  $s$ 
07 return  $-1$ 
```

- 运行时间和简单算法相同

# 启发式方法

- Boyer和Moore向逆向简单算法中增加了启发式规则，得到了  $O(n+m)$  算法, 但其更复杂
- Horspool建议仅使用出现启发式规则
  - 在不匹配之后，将  $T[s + m - 1]$  对齐到模式  $P[0..m-2]$  中的最右出现
  - 例：
    - $T = \text{"detective date"} , P = \text{"date"}$
    - $T = \text{"tea kettle"} , P = \text{"kettle"}$



# 偏移表

- 在预处理中, 计算大小为  $|\Sigma|$  的偏移表.

$$\text{shift}[w] = \begin{cases} m-1 - \max\{i < m-1 \mid P[i] = w\} & \text{if } w \text{ is in } P[0..m-2] \\ m & \text{otherwise} \end{cases}$$

- 例:  $P = \text{"kettle"}$ 
  - $\text{shift}[\mathbf{e}] = 4$ ,  $\text{shift}[\mathbf{l}] = 1$ ,  $\text{shift}[\mathbf{t}] = 2$ ,  $\text{shift}[\mathbf{k}] = 5$
- 例:  $P = \text{"pappar"}$ 
  - 其偏移表是什么?

# Boyer-Moore-Horspool 算法

**BMH-Search**(T, P)

```
01 // compute the shift table for P
01 for c  $\leftarrow$  0 to  $|\Sigma| - 1$ 
02     shift[c] = m          // default values
03 for k  $\leftarrow$  0 to m - 2
04     shift[P[k]] = m - 1 - k
05 // search
06 s  $\leftarrow$  0
07 while s  $\leq$  n - m do
08     j  $\leftarrow$  m - 1    // start from the end
09     // check if T[s..s+m-1] = P[0..m-1]
10     while T[s+j] = P[j] do
11         j  $\leftarrow$  j - 1
12         if j < 0 return s
13     s  $\leftarrow$  s + shift[T[s + m-1]]    // shift by last letter
14 return -1
```

# BMH 分析

- 最坏情况运行时间
  - 预处理:  $O(|\Sigma| + m)$
  - 搜索:  $O(nm)$ 
    - 何种输入达到此界?
  - 总计:  $O(nm)$
- 空间:  $O(|\Sigma|)$ 
  - 和  $m$  独立
- 在真实数据集集合上很快