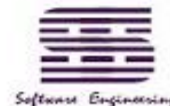


# 实用算法设计——查找

主讲：娄文启

[louwenqi@ustc.edu.cn](mailto:louwenqi@ustc.edu.cn)



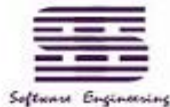
# 5 查找

**假定：待查找的数据已在内存中。**

**重点：**理解KMP算法和BM算法的思想。能利用KMP算法和BM算法解决实际的字串查找需求。

**难点：**理解KMP算法为何比朴素匹配算法快；BM算法为何比KMP算法快。

**基础：**C语言编程；线性结构的定义及基本操作的实现。



# 5 查找

## 5.1 基于Hash表的查找

## 5.2 蛮力查找（顺序查找）

## 5.3 基于有序表的二分查找

## 5.4 字符串的查找

## 5.5 基于树的查找



2024/11/27



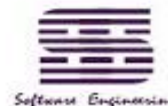
Software Engineering

- 字符串
- 模式匹配算法
  1. 朴素的模式匹配算法
  2. **KMP**模式匹配算法
  3. **BM**模式匹配算法

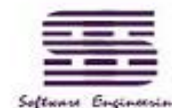


# 重难点

- 理解字符串是特殊的线性表，其特殊性在哪？
- 朴素的模式匹配算法为什么慢？
- **KMP**算法
  - 与朴素的模式匹配算法相比，**KMP**算法为什么快？
  - **KMP**算法中**Next**数组的设计初衷是？
  - 如何计算**Next**数组？如何利用**Next**数组来实现快速匹配？
  - **Next**数组存在什么问题，以致于后面提出了改进版的**NextVal**数组
  - 如何计算**NextVal**数组？如何利用**NextVal**数组来实现快速匹配？
- **BM**算法：
  - 理解两类启发式规则：“坏字符”和“好后缀”
  - 如何计算目标串上查找指针的跳跃举例**dist[i]**，以定位下轮匹配的右对齐的位置？
    - 如何计算**MatchJump**数组？
    - 如何计算**CharJump**数组



- 字符串：由零个或多个字符组成的有限序列。
  - 字符串  $s = "a_1a_2...a_n"$  ( $n \geq 0$ )
  - 串的长度：串中的字符数目  $n$
  - 空串：零个字符的串。即串的长度为0.
  - 空格串：只包含空格的串。
- 逻辑结构：可认为是特殊的线性表。
  - 串中的元素只能是字符类型。
  - 字符串：关注子串的操作：查找/插入/删除/替换子串。
    - 线性表：关注单个元素的操作
- 存储结构：推荐顺序存储方式。
  - 例如，存储串 "good"，需要5个字节，但是其长度为4。



# 模式匹配算法

- 字符串查找场景：
  - 例如，在文本“**fffffab cfe defe**”中查找字符串“**ff**”
  - 本质：在字符线性表（主串或目标字符串T）中**查找**匹配的子表（子串或模式字符串P）。
- 字符串查找（或模式匹配）算法：
  1. 朴素的模式匹配算法
  2. KMP算法（Knuth-Morris-Pratt)
  3. BM算法（Boyer-Moore)



# 1. 朴素的模式匹配算法

## 示例1:

假设我们要从下面的主串 **T**="goodgoogle" 中，找到 **P**="google" 这个子串的位置。我们通常需要下面的步骤。

1. 主串 **T** 第一位开始，**T** 与 **P** 前三个字母都匹配成功，但 **T** 第四个字母是 **d** 而 **P** 的是 **g**。第一位匹配失败。如图 5-6-1 所示，其中竖直连线表示相等，闪电状弯折连线表示不等。

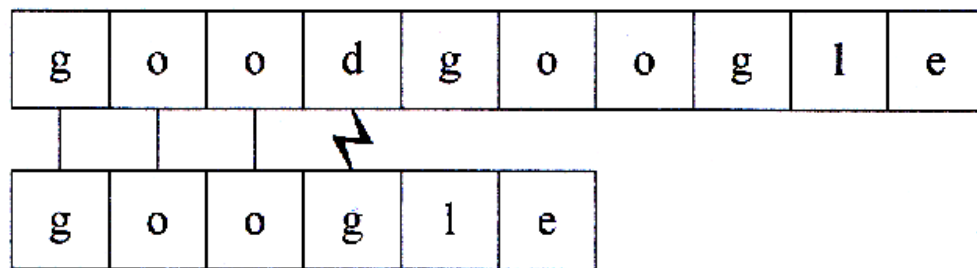


图 5-6-1



2. 主串 **T** 第二位开始，主串 **T** 首字母是 o，要匹配的 **P** 首字母是 g，匹配失败，如图 5-6-2 所示。

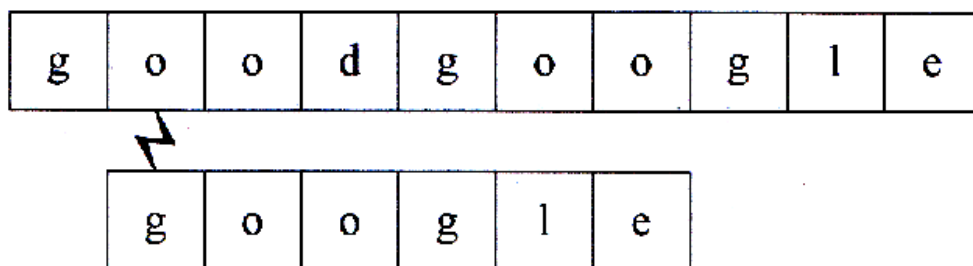


图 5-6-2

3. 主串 **T** 第三位开始，主串 **T** 首字母是 o，要匹配的 **P** 首字母是 g，匹配失败，如图 5-6-3 所示。

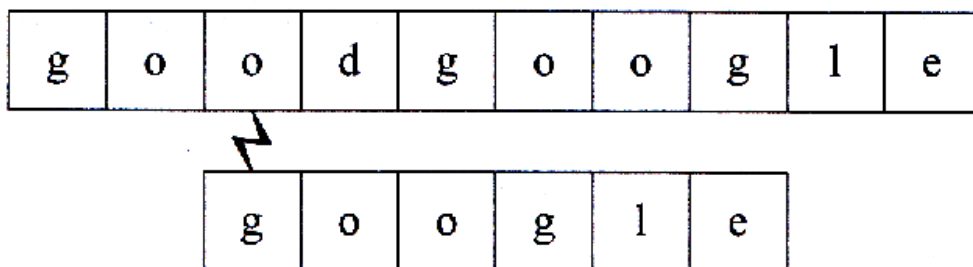


图 5-6-3



4. 主串 **T** 第四位开始，主串 **T** 首字母是 **d**，要匹配的 **P** 首字母是 **g**，匹配失败，如图 5-6-4 所示。

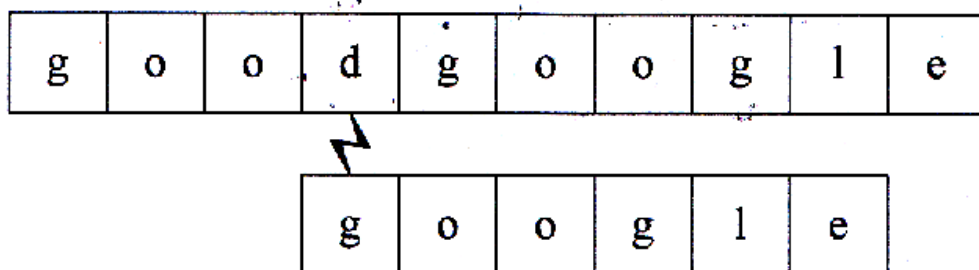


图 5-6-4

5. 主串 **T** 第五位开始，**T** 与 **P**，6 个字母全匹配，匹配成功，如图 5-6-5 所示。

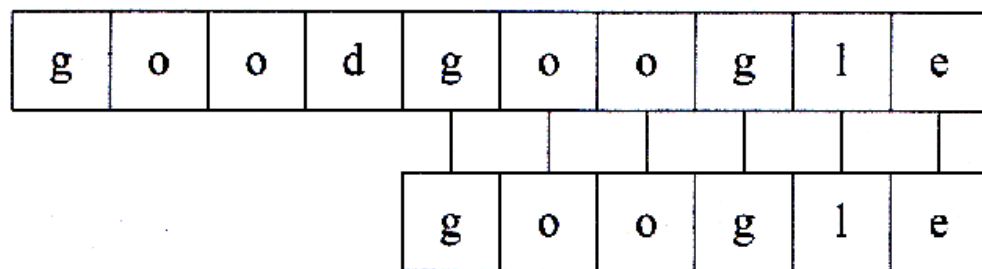


图 5-6-5

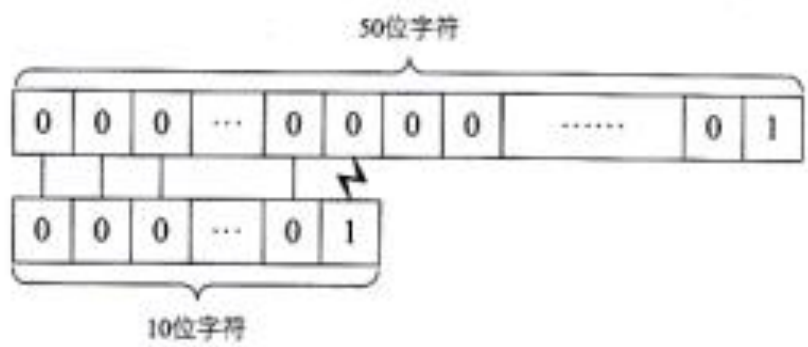


假设目标串T长度为n，模式串P长度为m。

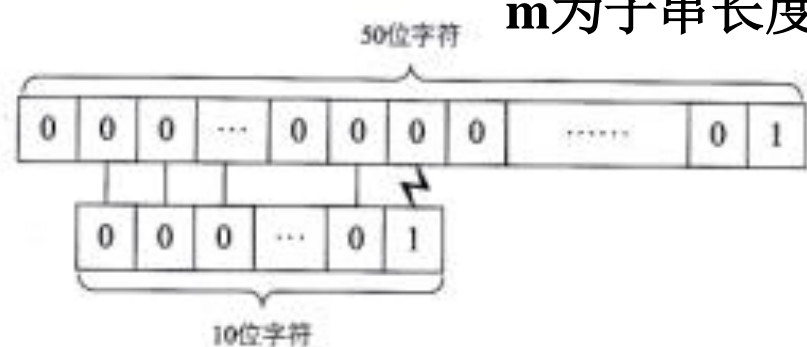
朴素的模式匹配算法的最坏情况下近似时间复杂度为 $O((n-m+1)*m)$ 。

示例2:

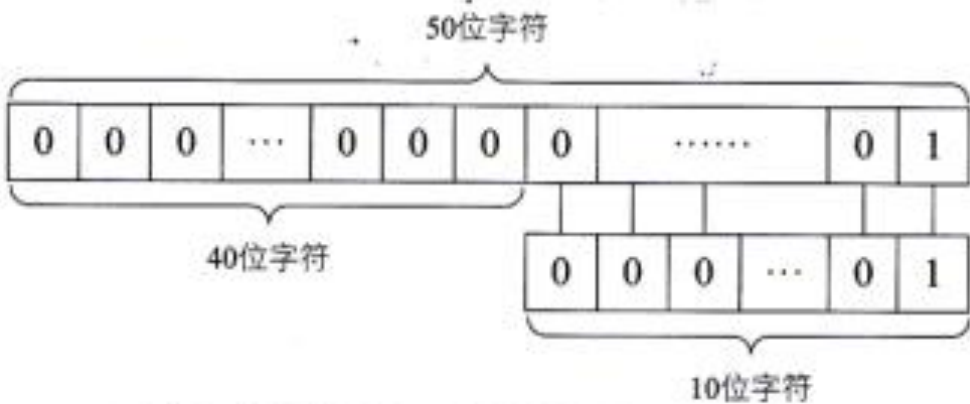
(假设n为主串长度，  
m为子串长度)



T在第一位置判断了10次发现字符串不匹配



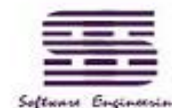
T在第二位置再次判断了10次发现字符串不匹配



T在第41个位置判断了10次发现字符串终于匹配成功。  
期间进行了  $(50-10+1) \times 10$  次判断操作



- 假设目标串**T**长度为**n**，模式串**P**长度为**m**。  
假定**T**中确实存在子串**P**（即假定模式匹配成功），则：
  - 最好的情况下的近似时间复杂度为 $O(m)$ 
    - 比如，在"googlegood"中查找"google".
  - 若每次不成功的匹配都发生在串**P**的首字符处，  
则平均情况下的近似时间复杂度为 $O(n+m)$
  - 最坏的情况下的近似时间复杂度为 $O((n-m+1)*m)$ .
    - 比如，在“000...00001"中查找“0001".



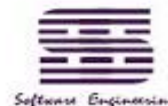
# 朴素的模式匹配算法 -> KMP, BM

- 朴素的模式匹配算法：
  - 需进行回溯
- 字符串查找算法中，最著名的两个算法：**KMP算法（Knuth-Morris-Pratt）**和**BM算法（Boyer-Moore）**
  - 它们都是精确字符串匹配算法（区别于模糊匹配）。
  - 目标字符串中无需进行回溯。（都比朴素的模式匹配算法快）
  - 模式字符串的移动方向：从目标字符串的第一个字符开始，朝目标字符串的尾部方向移动搜索匹配子串。
  - 每轮匹配时，字符匹配的方向：
    - **KMP算法**：采用从左向右进行字符的匹配比较。
    - **BM算法**：采用从右向左进行字符的匹配比较。

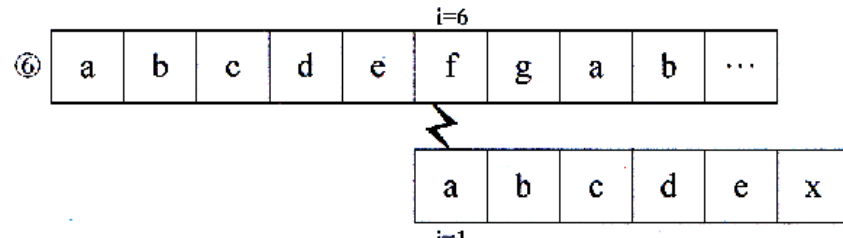
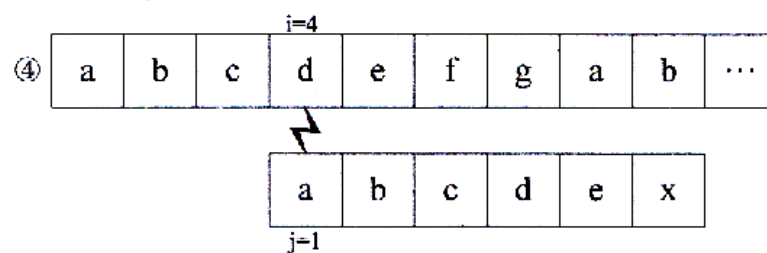
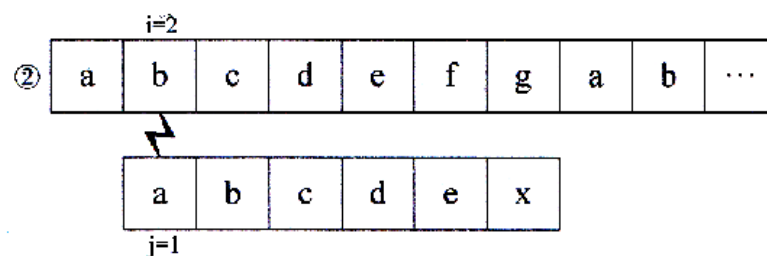
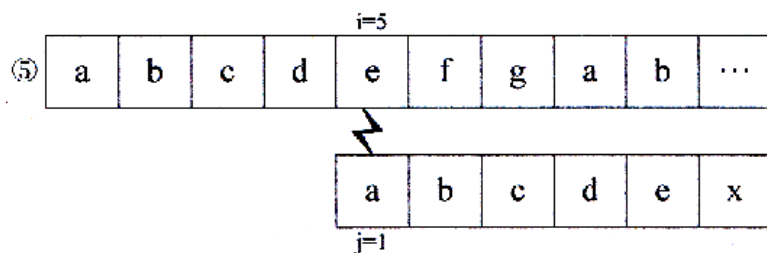
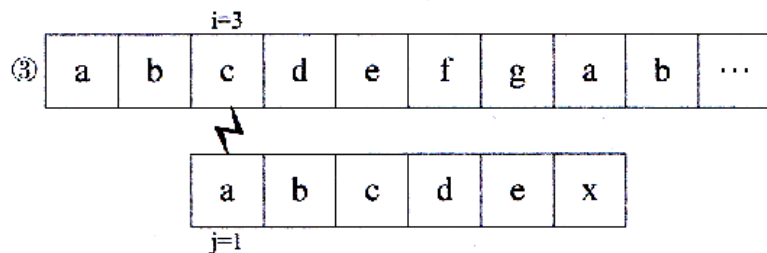
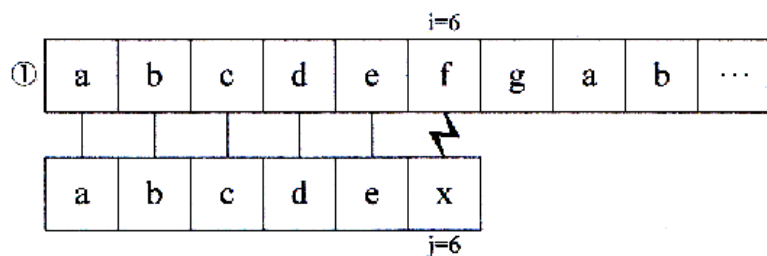


## 2. KMP算法

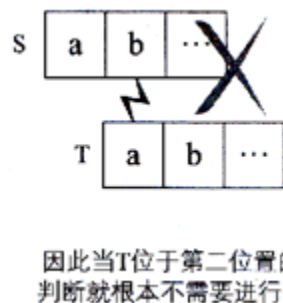
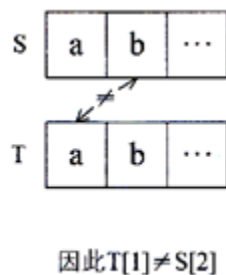
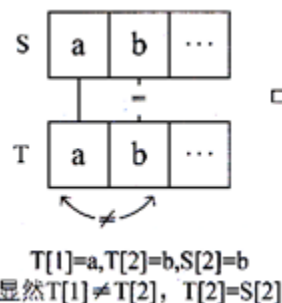
- **KMP算法的设计初衷：**
  - 希望能在匹配过程中，目标串**T**中的字符匹配是一直向右前进的，而不会出现向左的回溯
  - 理论上可以避免吗？
  - 如何避免？（利用模式串**P**自身的重复模式）



# 例1：模式串P中无重复模式

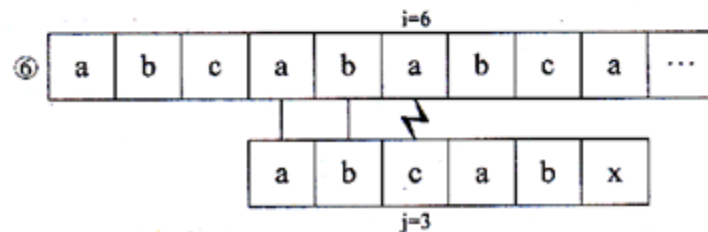
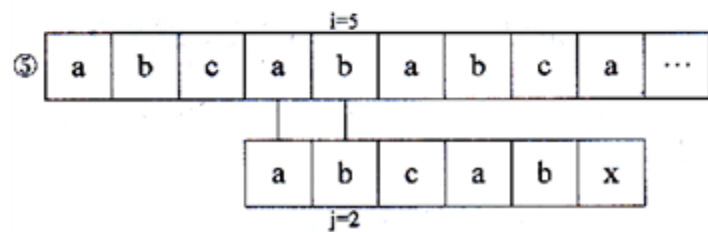
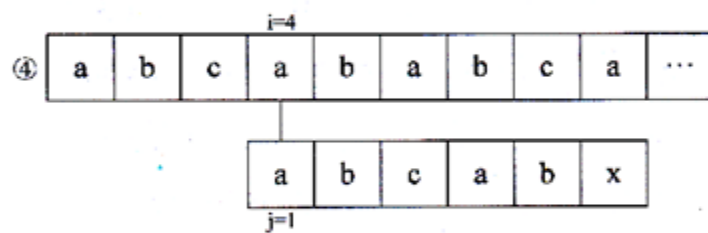
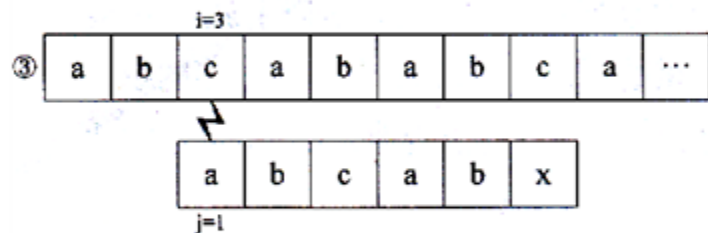
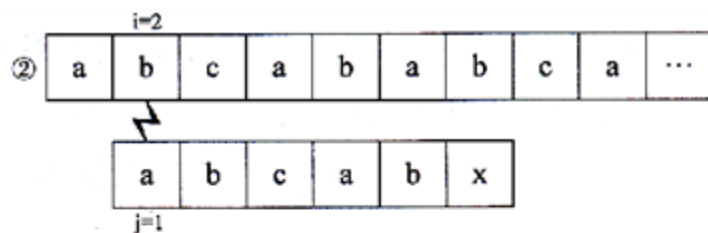
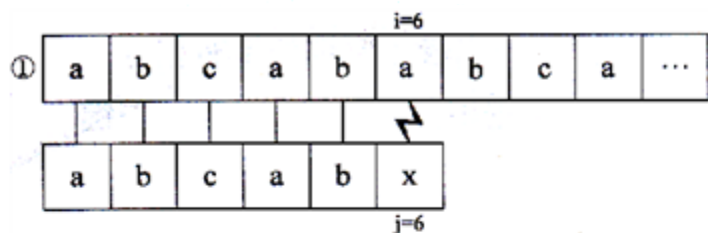


②③④⑤的判断都是多余

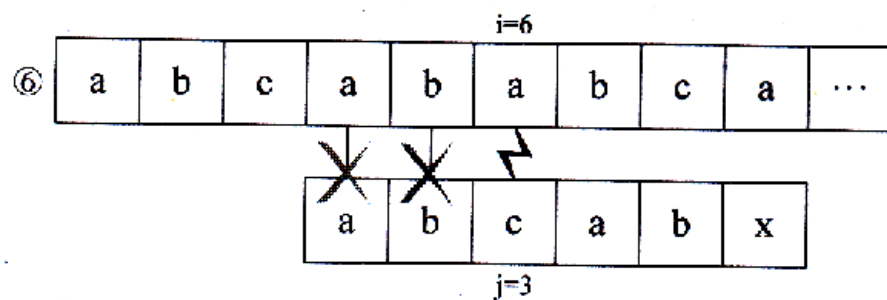
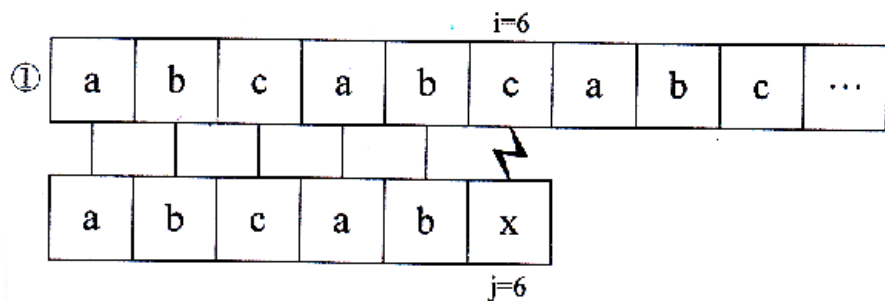


## 例2：模式串P中有重复模式

假设  $T = \text{"abcabcabc"}$ ,  $P = \text{"abcabx"}$ 。

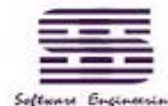


步骤是②③④⑤都是多余的！⑥中的前2次字母的匹配也是多余的！

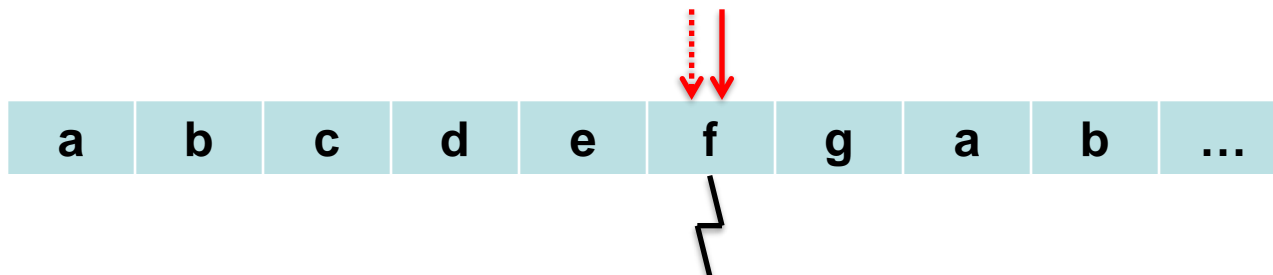




- 如何避免目标串中字符匹配的回溯问题？  
（利用模式串**P**自身的重复模式）
  - 假定目标串为**T**，模式串为**P**。若**T[i]≠P[j]**导致当前轮的匹配失败，则下一轮匹配时应该比较**T**和**P**中的哪两个字符？
    - 每轮中，**T**与**P**都是从左至右逐字符比较。
  - 关注：
    - 匹配失败，是否发生在**P**的第一个字符处？
    - **P**中是否有重复模式？



## • Case 1



第一轮:

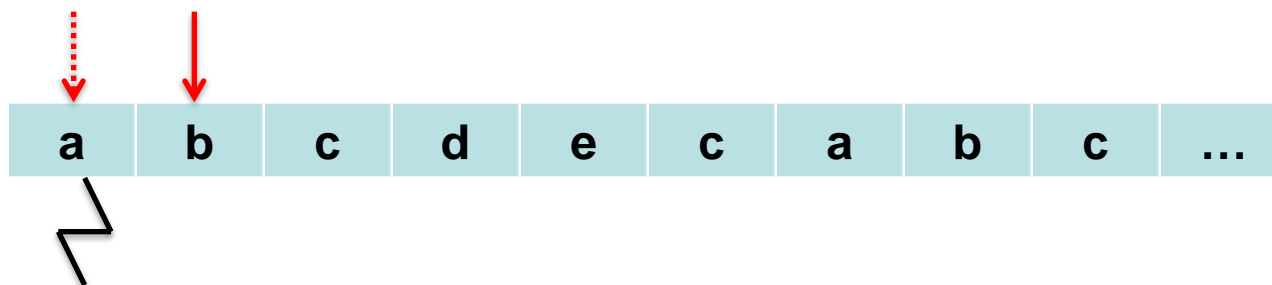


第二轮:



模式串中无重复模式

## • Case 2



第一轮:



第二轮:



模式串中第一个字符不匹配



## • Case 3

第一轮:



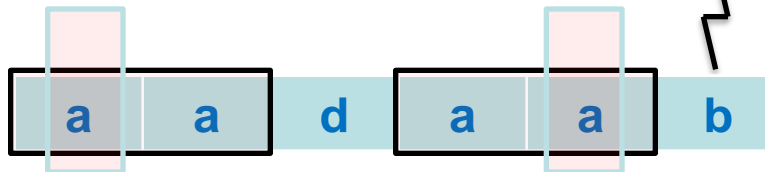
第二轮:



模式串中存在重复模式

## • Case 4

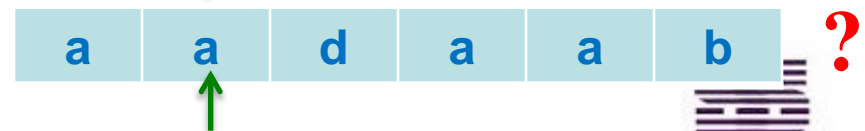
第一轮:



第二轮:



?



?

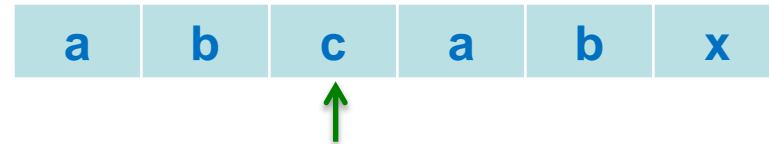


## • Case 3

第一轮:



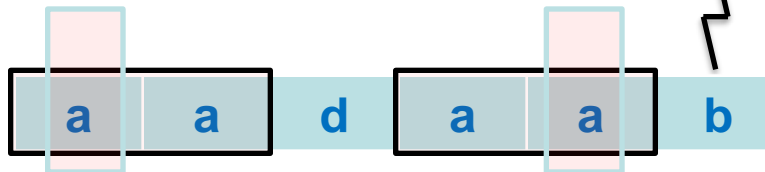
第二轮:



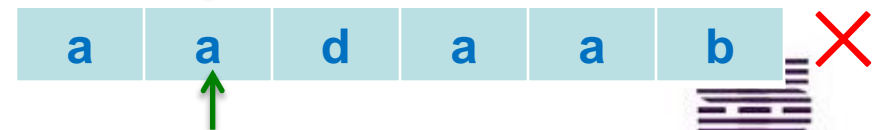
模式串中存在重复模式

## • Case 4

第一轮:



第二轮:

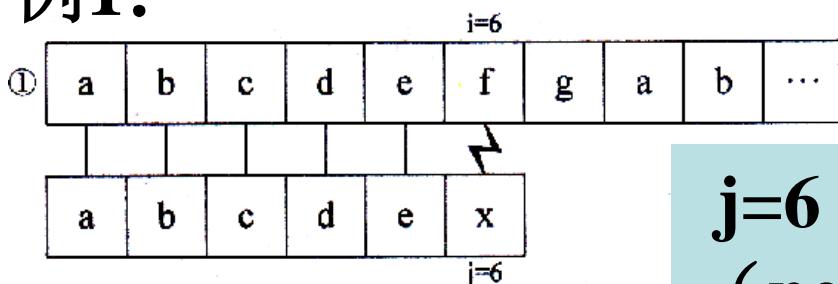


- 如何避免目标串中字符匹配的回溯问题？（利用模式串**P**自身的重复模式）
  - 假定目标串为**T**，模式串为**P**。若 **$T[i] \neq P[j]$** 导致当前轮的匹配失败，则下一轮匹配时应该比较**T**和**P**中的哪两个字符？
    - 每轮中，**T**与**P**都是从左至右逐字符比较。
  - 案例分析总结：
    - 若当前轮匹配在进行第一个字符比较时就失败，那么下一轮应该也是比较 **$T[i+1]$** 和 **$P[1]$**
    - 若**P**中在当前轮成功匹配的子串的后缀与子串的前缀无重复模式，那么下一轮应该也是比较 **$T[i]$** 和 **$P[1]$**
    - 若**P**中在当前轮成功匹配的子串的后缀与子串的前缀有重复模式，那么下一轮应该也是比较 **$T[i]$** 和 **$P[\text{next}[j]]$** 
      - **$\text{next}[j] = ?$**

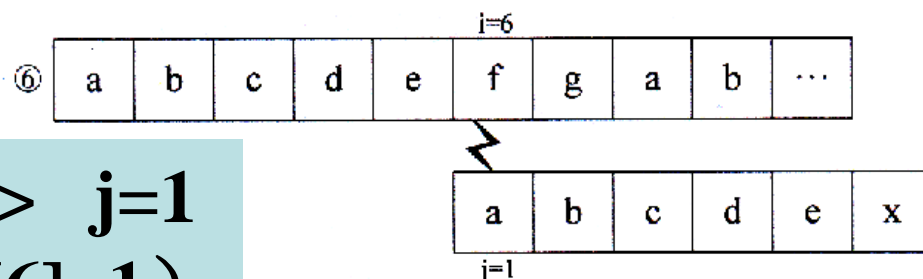
# next数组

**注意：**在KMP算法中，模式串P的下标从1开始，而不是从0开始。  
为什么呢？

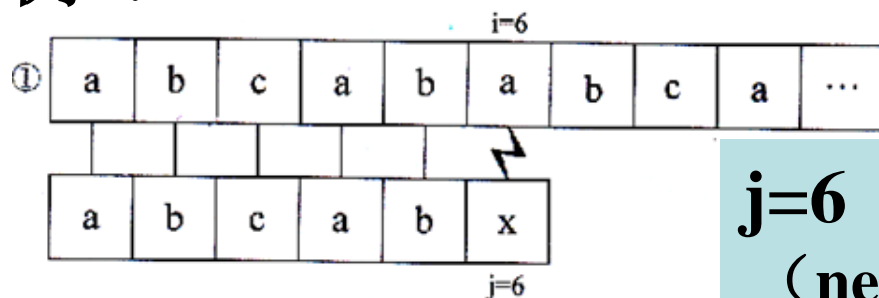
例1:



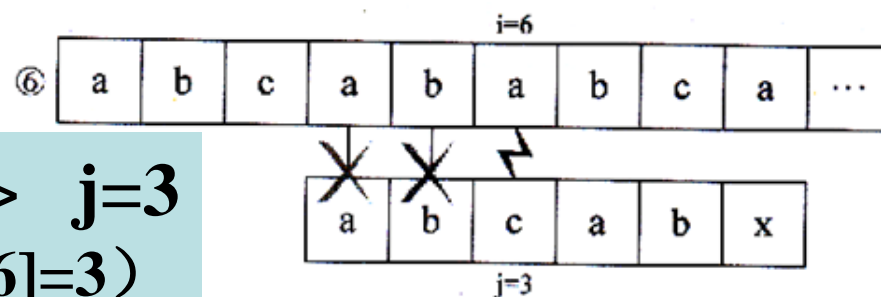
**j=6 => j=1**  
(next[6]=1)



例2:



**j=6 => j=3**  
(next[6]=3)



我们把 T 串各个位置的 j 值的变化定义为一个数组 next，那么 next 的长度就是 T 串的长度。于是我们可以得到下面的函数定义：

$$next[j] = \begin{cases} 0, & \text{当 } j=1 \text{ 时} \\ \text{Max} \{ k \mid 1 < k < j, \text{ 且 } 'p_1 \dots p_{k-1}' = 'p_{j-k+1} \dots p_{j-1}' \} & \text{当此集合不空时 (隐含 } j > 2) \\ 1, & \text{其他情况} \end{cases}$$

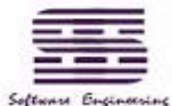
**Next[j]表征：**下一轮匹配时比较P与T的哪两个字符？

# KMP算法的主要思想

- 将模式串P自身的重复规律保存到next数组中

$$next[j] = \begin{cases} 0, & \text{当 } j=1 \text{ 时} \\ \text{Max} \{ k \mid 1 < k < j, \text{ 且 } 'p_1 \cdots p_{k-1}' = 'p_{j-k+1} \cdots p_{j-1}' \} & \text{当此集合不空时} \\ 1, & \text{其他情况} \end{cases}$$

- 匹配过程：若某轮匹配失败，则利用next数组分别计算下一轮匹配时目标串和模式串的开始位置
  - 若是 $T[i] \neq P[j]$ 导致当前轮的匹配失败，则按照下列规则开始下一轮匹配：
    - 若 $next[j] \neq 0$ ，则将 $T[i..]$ 与 $P[next[j]..]$ 匹配；
    - 若 $next[j] == 0$ ，则将 $T[(i+1)..]$ 与 $P[1..]$ 匹配。



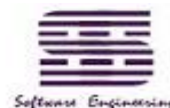
# next数组的计算示例

1.  $P = \text{"abcdex"}$  (如表 5-7-1 所示)

表 5-7-1

j	123456
模式串	abcdex
next[j]	011111

- 1) 当  $j=1$  时,  $\text{next}[1]=0$ ;
- 2) 当  $j=2$  时,  $j$  由 1 到  $j-1$  就只有字符 “a”, 属于其他情况  $\text{next}[2]=1$ ;
- 3) 当  $j=3$  时,  $j$  由 1 到  $j-1$  串是 “ab”, 显然 “a” 与 “b” 不相等, 属其他情况,  $\text{next}[3]=1$ ;
- 4) 以后同理, 所以最终此 T 串的  $\text{next}[j]$  为 011111。





## 2. $P = \text{"abcbx"}$ (如表 5-7-2 所示)

表 5-7-2

j	123456
模式串	abcbx
next[j]	011123

- 1) 当  $j=1$  时,  $\text{next}[1]=0$ ;
- 2) 当  $j=2$  时, 同上例说明,  $\text{next}[2]=1$ ;
- 3) 当  $j=3$  时, 同上,  $\text{next}[3]=1$ ;
- 4) 当  $j=4$  时, 同上,  $\text{next}[4]=1$ ;
- 5) 当  $j=5$  时, 此时  $j$  由 1 到  $j-1$  的串是 "abca", 前缀字符 "a" 与后缀字符 "a" 相等 (前缀用下划线表示, 后缀用斜体表示), 因此可推算出  $k$  值为 2 (由 ' $p_1 \cdots p_{k-1}$ ' = ' $p_{j-k+1} \cdots p_{j-1}$ ', 得到  $p_1=p_4$ ) 因此  $\text{next}[5]=2$ ;
- 6) 当  $j=6$  时,  $j$  由 1 到  $j-1$  的串是 "abcab", 由于前缀字符 "ab" 与后缀 "ab" 相等, 所以  $\text{next}[6]=3$ 。

### 3. P="ababaaaba" (如表 5-7-3 所示)

表 5-7-3

j	123456789
模式串	ababaaaba
next[j]	011234223

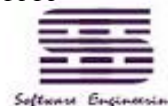
- 1) 当  $j=1$  时,  $\text{next}[1]=0$ ;
- 2) 当  $j=2$  时, 同上  $\text{next}[2]=1$ ;
- 3) 当  $j=3$  时, 同上  $\text{next}[3]=1$ ;
- 4) 当  $j=4$  时,  $j$  由 1 到  $j-1$  的串是 "aba", 前缀字符 "a" 与后缀字符 "a" 相等,  $\text{next}[4]=2$ ;
- 5) 当  $j=5$  时,  $j$  由 1 到  $j-1$  的串是 "abab", 由于前缀字符 "ab" 与后缀 "ab" 相等, 所以  $\text{next}[5]=3$ ;
- 6) 当  $j=6$  时,  $j$  由 1 到  $j-1$  的串是 "ababa", 由于前缀字符 "aba" 与后缀 "aba" 相等, 所以  $\text{next}[6]=4$ ;
- 7) 当  $j=7$  时,  $j$  由 1 到  $j-1$  的串是 "ababaa", 由于前缀字符 "ab" 与后缀 "aa" 并不相等, 只有 "a" 相等, 所以  $\text{next}[7]=2$ ;
- 8) 当  $j=8$  时,  $j$  由 1 到  $j-1$  的串是 "ababaaa", 只有 "a" 相等, 所以  $\text{next}[8]=2$ ;
- 9) 当  $j=9$  时,  $j$  由 1 到  $j-1$  的串是 "ababaaab", 由于前缀字符 "ab" 与后缀 "ab" 相等, 所以  $\text{next}[9]=3$ 。

4.  $P = \text{"aaaaaaaaab"}$  (如表 5-7-4 所示)

表 5-7-4

j	123456789
模式串	aaaaaaaaab
next[j]	012345678

- 1) 当  $j=1$  时,  $\text{next}[1]=0$ ;
- 2) 当  $j=2$  时, 同上  $\text{next}[2]=1$ ;
- 3) 当  $j=3$  时,  $j$  由 1 到  $j-1$  的串是 "aa", 前缀字符 "a" 与后缀字符 "a" 相等,  $\text{next}[3]=2$ ;
- 4) 当  $j=4$  时,  $j$  由 1 到  $j-1$  的串是 "aaa", 由于前缀字符 "aa" 与后缀 "aa" 相等, 所以  $\text{next}[4]=3$ ;
- 5) .....
- 6) 当  $j=9$  时,  $j$  由 1 到  $j-1$  的串是 "aaaaaaaa", 由于前缀 "aaaaaaaa" 与后缀 "aaaaaaaa" 相等, 所以  $\text{next}[9]=8$ .



# next数组的应用

- 例：在目标串 “aaaaaaaaacaaaaaaaaaab” 中查找模式串 “aaaaaaaaab” 的位置。
- **Step1:** 计算模式串的next数组：

j	123456789
模式串	aaaaaaaaab
next[j]	012345678

- **Step2:** 匹配过程。（9轮）



# KMP算法的实现

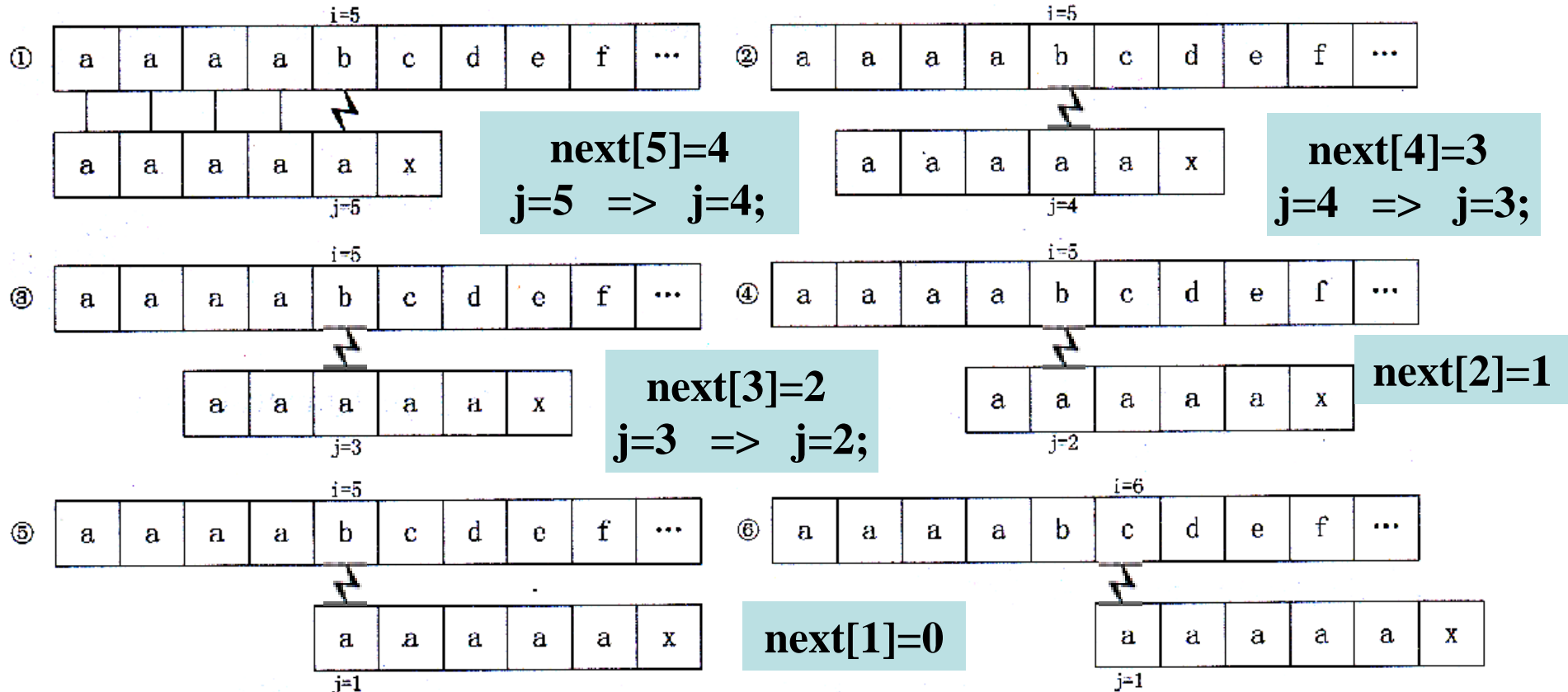
```
class Solution {
public:
    int strStr(string s, string p) {
        int n,m;
        n = s.size();
        m = p.size();
        vector<int> next;
        int i=0,j=-1;
        next.push_back(-1);
        while(i < m){
            if(j==-1 or p[i] == p[j]){
                i++;
                j++;
                next.push_back(j);
            }else {
                j = next[j];
            }
        }
    }
};
```

```
i=0;j=0;
while( i < n and j < m){
    if(j==-1 || s[i] == p[j]){
        i++;
        j++;
    }else{ j = next[j]; }
}

return (j==m) ? i-m : -1;
// return 0 ;
```



# Q: next数组还有改进空间吗?

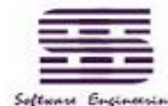


步骤②③④⑤都是多余的!

我们发现，当中的②③④⑤步骤，其实是多余的判断。由于 T 串的第二、三、四、五位置的字符都与首位的“a”相等，那么可以用首位 next[1]的值去取代与它相等的字符后续 next[j]的值，这是个很好的办法。因此我们对求 next 函数进行了改良。

# KMP算法的改进

- 什么情况下有改进的空间？
  - 假设 $T[i] \neq P[j]$ 导致失配。若 $P[j] == P[\text{next}[j]]$ ，此时若向右移动模式串 $P$ ，将 $T[i]$ 与 $P[\text{next}[j]]$ 对齐进行比较必然是无意义的，因为此时 $T[i]$ 必定 $\neq P[\text{next}[j]]$ 。
- 如何改进？用 $\text{nextval}$ 数组代替 $\text{next}$ 数组。
  - $\text{nextval}[1]=0$ ;
  - $\text{for}(j>1; j \leq n; j++)$ 
    - 若 $P[j] == P[\text{next}[j]]$ ，则 $\text{nextval}[j] = \text{nextval}[\text{next}[j]]$ ;
    - 若 $P[j] \neq P[\text{next}[j]]$ ，则 $\text{nextval}[j] = \text{next}[j]$ ;
- 若某轮匹配失败，则利用 $\text{nextval}$ 数组计算下一轮匹配时的目标串和模式串的开始位置（类似 $\text{next}$ 数组的应用）



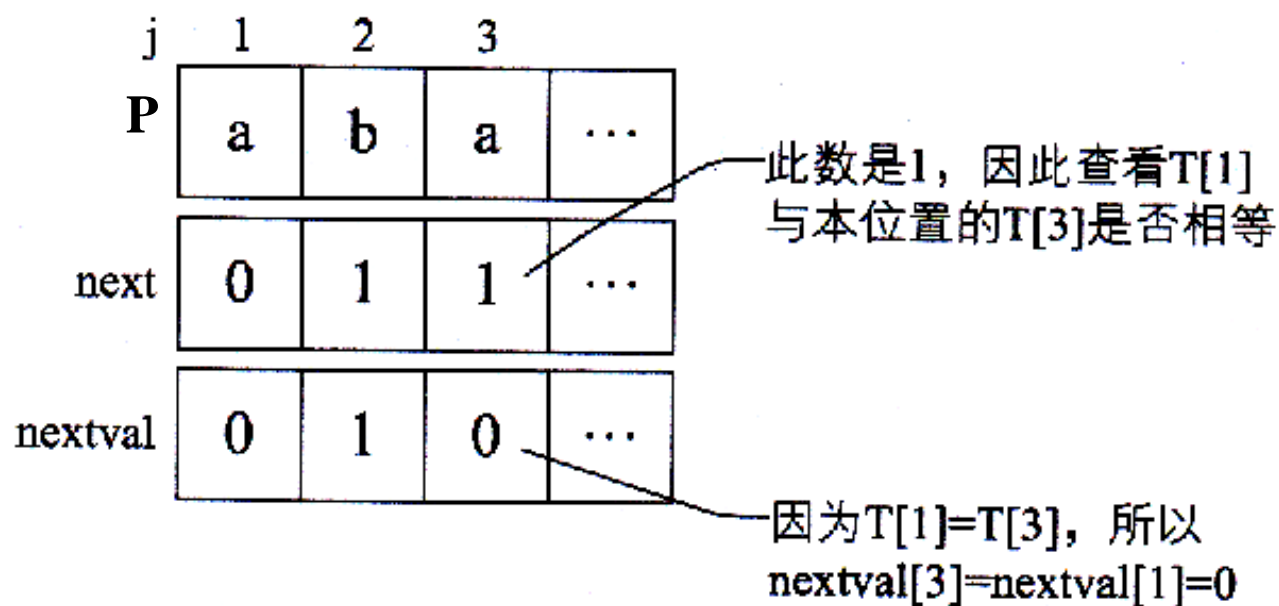


图 5-7-7

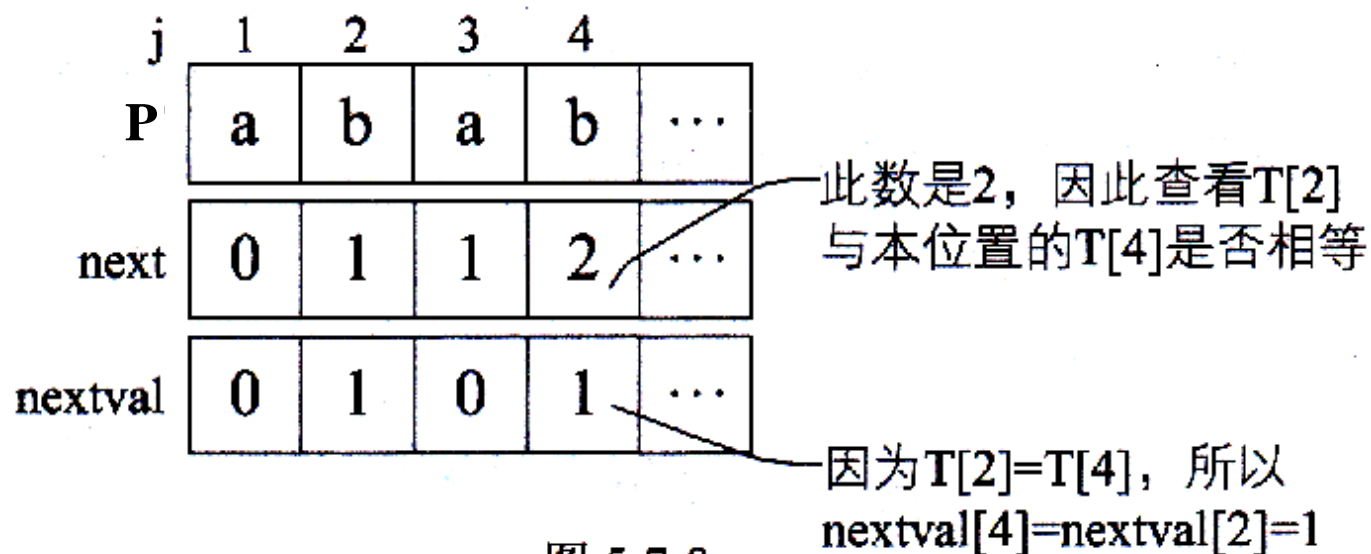


图 5-7-8





# nextval数组的计算示例

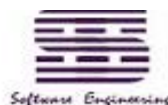
1. P="ababaaaba" (如表 5-7-5 所示)

表 5-7-5

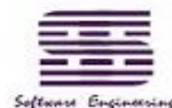
j	123456789
模式串	ababaaaba
next[j]	011234223
nextval[j]	010104210

先算出 next 数组的值分别为 001234223，然后再分别判断。

- 1) 当  $j=1$  时,  $\text{nextval}[1]=0$ ;
- 2) 当  $j=2$  时, 因第二位字符“b”的 next 值是 1, 而第一位就是“a”, 它们不相等, 所以  $\text{nextval}[2]=\text{next}[2]=1$ , 维持原值。
- 3) 当  $j=3$  时, 因为第三位字符“a”的 next 值为 1, 所以与第一位的“a”比较得知它们相等, 所以  $\text{nextval}[3]=\text{nextval}[1]=0$ ; 如图 5-7-7 所示。
- 4) 当  $j=4$  时, 第四位的字符“b”next 值为 2, 所以与第二位的“b”相比较得到结果是相等, 因此  $\text{nextval}[4]=\text{nextval}[2]=1$ ; 如图 5-7-8 所示。



- 5) 当  $j=5$  时,  $next$  值为 3, 第五个字符 “a” 与第三个字符 “a” 相等, 因此  $nextval[5]=nextval[3]=0$ ;
- 6) 当  $j=6$  时,  $next$  值为 4, 第六个字符 “a” 与第四个字符 “b” 不相等, 因此  $nextval[6]=4$ ;
- 7) 当  $j=7$  时,  $next$  值为 2, 第七个字符 “a” 与第二个字符 “b” 不相等, 因此  $nextval[7]=2$ ;
- 8) 当  $j=8$  时,  $next$  值为 2, 第八个字符 “b” 与第二个字符 “b” 相等, 因此  $nextval[8]=nextval[2]=1$ ;
- 9) 当  $j=9$  时,  $next$  值为 3, 第九个字符 “a” 与第三个字符 “a” 相等, 因此  $nextval[9]=nextval[3]=1$ 。



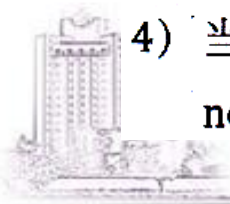
## 2. P="aaaaaaaaab" (如表 5-7-6)

表 5-7-6

j	123456789
模式串 T	aaaaaaaaab
next[j]	012345678
nextval[j]	000000008

先算出 next 数组的值分别为 012345678，然后再分别判断。

- 1) 当 j=1 时, nextval[1]=0;
- 2) 当 j=2 时, next 值为 1, 第二个字符与第一个字符相等, 所以 nextval[2]=nextval[1]=0;
- 3) 同样的道理, 其后都为 0……;
- 4) 当 j=9 时, next 值为 8, 第九个字符“b”与第八个字符“a”不相等, 所以 nextval[9]=8。



# 例题：计算nextval数组

• 1.

j	123456
模式串	abcbabx
next[j]	011123

• 2.

j	123456
模式串	abcdex
next[j]	011111

• 3.

j	123456789
模式串	ababaaaba
next[j]	011234223

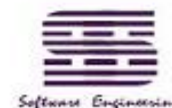


# nextval数组的应用

- 例：在目标串 “aaaaaaaaacaaaaaaaaaab” 中查找模式串 “aaaaaaaaab” 的位置。
- **Step1:** 计算模式串的next数组：

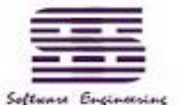
j	123456789
模式串	aaaaaaaaab
next[j]	012345678
nextval[j]	000000008

- **Step2:** 匹配过程。（2轮）



# KMP算法的性能 vs. 朴素算法

- 近似时间复杂度为 $O(n+m)$ , 其中 $O(n)$ 表示比较的时间,  $O(m)$ 表示计算next数组的时间.
  - 目标串中不存在回溯情况
  - 目标串中的每个字符, 会比较1~2次;
- **注意:** KMP算法中, 字符串的下标从1开始, 而不是从0开始
- 仅当模式串与目标串之间存在许多“部分匹配”情况下, KMP算法才比朴素模式匹配算法更具优势。
- 若每轮中模式串与目标串之间的不匹配都发生在模式串的第一个字符处, 则KMP算法会退化到朴素模式匹配算法。
  - 因为 $nextval[1]=0$ ;



# 3. Boyer-Moore查找

- 是字符串查找算法中最著名的两个算法之一
- **BM算法**:
  - 是一种**精确字符串匹配**算法（区别于模糊匹配）。
  - 每轮匹配过程中，字符的匹配方向：采用**从右向左**进行字符的匹配比较。（与**KMP**算法的主要区别）
    - 若某轮匹配失败，则移动模式串，与目标串的下一轮开始匹配位置进行右对齐，然后开始当前轮的从右至左的逐字符匹配。
    - 故，**BM**算法中的关键问题是，如何确定目标串中的下一轮匹配的起始位置？即，如何确定目标串中查找指针的移动距离？
  - 采用**启发式方法**：无需检查目标字符串中的所有字符
    - 利用**P**中的重复模式和 **T**中的失配字符



# Boyer-Moore示例

假定字符串为"HERE IS A SIMPLE EXAMPLE", 搜索词为"EXAMPLE".

1

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

"S"与"E"不匹配。这时, "S"就被称为"坏字符", 即不匹配的字符  
"S"不包含在搜索词"EXAMPLE"之中,

2

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

模式串移动了7位



# Boyer-Moore示例

2

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

依然从尾部开始比较，发现"P"与"E"不匹配，所以"P"是"坏字符"。但，"P"包含在搜索词"EXAMPLE"之中。所以，将搜索词后移两位，两个"P"对齐。

3

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

模式串移动了2位



# Boyer-Moore示例

坏字符规则：当文本串中的某个字符跟模式串的某个字符不匹配时，我们称文本串中的这个失配字符为坏字符。

此时模式串需要向右移动，

移动的位数 = 坏字符在模式串中的位置 - 坏字符在模式串中最右出现的位置。

此外，如果"坏字符"不包含在模式串之中，则最右出现位置为-1。

后移位数 = 坏字符的位置 - 搜索词中的上一次出现位置



# Boyer-Moore示例

3

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

4

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

仍然从最右侧开始  
比较

5

HERE IS A SIMPLE EXAMPLE  
EXAMPLE



# Boyer-Moore示例

7

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

8

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

不匹配出现

坏字符规则，模式串向右3位（2-(-1)）

如考虑已经匹配的“MPLE”后缀，是否有更好的方式

“MPLE”、“PLE”，“LE”，“E”



2024/11/27



44

# Boyer-Moore示例

8

HERE IS A SIMPLE EXAMPLE

EXAMPLE

9

HERE IS A SIMPLE EXAMPLE

EXAMPLE

好后缀"的位置以最后一个字符为准。

若"好后缀"在P只出现一次，则它上一次出现位置为 -1

所有的"好后缀"（MPLE、PLE、LE、E）之中，只有"E"在  
"EXAMPLE"还出现在头部，所以后移  $6 - 0 = 6$  位。

选择最小的移动距离（而非7位）



# Boyer-Moore示例

8

HERE IS A SIMPLE EXAMPLE

EXAMPLE

9

HERE IS A SIMPLE EXAMPLE

EXAMPLE

好后缀规则：当字符失配时，

模式串后移位数

= 好后缀在模式串中的位置 - 好后缀在模式串上一次出现的位置

如果好后缀在模式串中没有再次出现，则为-1。



# Boyer-Moore示例

9

HERE IS A SIMPLE EXAMPLE  
EXAMPLE

10

HERE IS A SIMPLE EXAMPLE  
EXAMPLE



- **BM算法中的关键问题：**当匹配失败时，如何确定目标串中的下一轮匹配的起始位置？即，如何决定目标串中指针向右跳跃的距离**dist**？
  - 假设出现失配时， $T[i] \neq P[k]$ 。则此时坏字符为 $x(=T[i])$ ，好后缀 $P'=P[(k+1) \dots (\text{len}(P)-1)]$ （好后缀，是已匹配的部分字符串）
  - 如何确定目标串中查找指针的移动距离**dist[i]**？
  - 采用**2种启发式方法**：无需检查目标串中的所有字符即可查找是否存在匹配子串。
- 启发式方法#1：跳过字符（“坏字符”规则）
  - $\text{CharJump}[x]$ ：依据T中的坏字符x，计算T中查找指针的跳跃距离。
- 启发式方法#2：重复模式（“好后缀”规则）
  - $\text{MatchJump}[k]$ ：依据P中的失配位置k，计算T中查找指针的跳跃距离。
- 例1：在目标串“oaks from acorns grow”中查找模式串“corn”的位置。（模式串中无重复模式）
- 例2：在目标串“he went door to door”中查找模式串“door to door”的位置。（模式串具有重复模式）





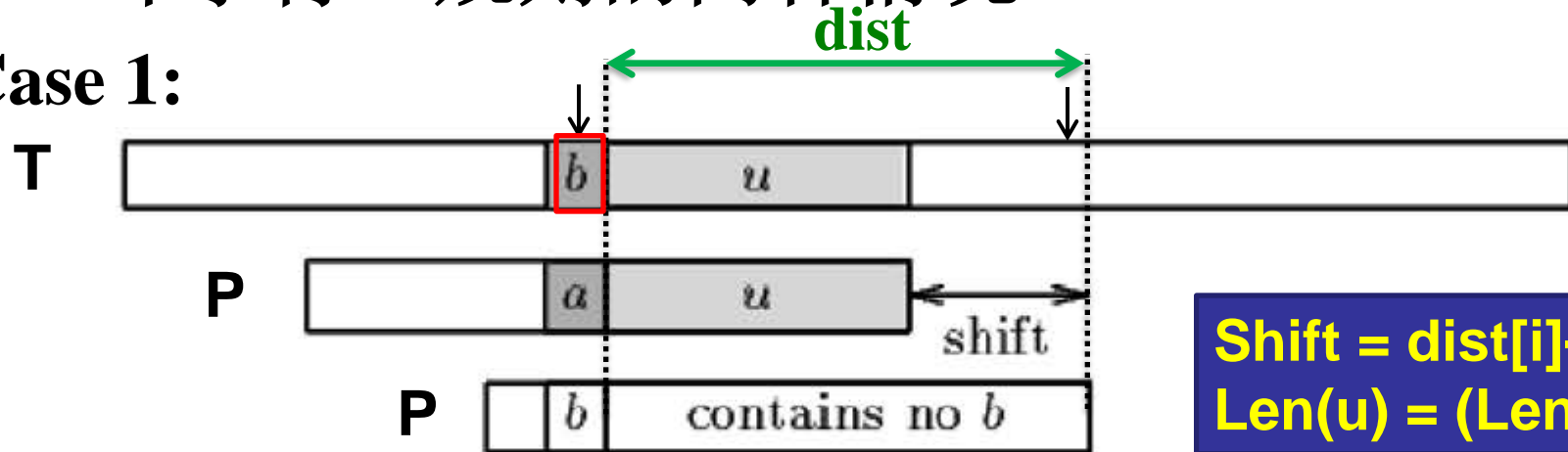
# BM算法的基本思想

- 设目标串T，模式串为P。
- **Step1:** 对于模式串P，计算CharJump[x]和MatchJump[k]。
- **Step2:** 将T与P的第一个字符对齐。
- **Step3:** T与P进行从右向左的逐字符比较，直至找到一个不匹配字符或者P中所有字符都匹配成功。
- **Step4:** 若出现失配，即存在 $T[i] \neq P[k]$ ，此时坏字符 $x = T[i]$ ，好后缀 $P' = P[(k+1) \dots (\text{len}(P)-1)]$ 。按如下规则计算目标串T中查找指针向右移动dist[i]:
  - 若此时T与P已有部分字符匹配（即存在“好后缀”）时，BM算法将采用2种启发式方法（即坏字符规则 和好后缀规则），计算 $\text{dist}[i] = \max(\text{CharJump}[x], \text{MatchJump}[k])$ 。
  - 若不存在“好后缀”，则必定是在模式串P的最后一个字符处出现失配。此时应采用启发式方法#1：跳过字符规则（“坏字符”规则），计算设置 $\text{dist}[i] = \text{CharJump}[x]$ 。
- **Step5:** 若 $(i + \text{dist}[i]) \leq \text{Len}(T) - 1$ ，则移动模式字符串P，使之与 $T[i + \text{dist}[i]]$ 右对齐，重复Step3；否则，认为T中不存在与P匹配的子串，返回匹配失败。

$$\text{dist}[i] = \text{CharJump}['b']$$

# “坏字符”规则的两种情况：

Case 1:



$$\text{Shift} = \text{dist}[i] - \text{Len}(u),$$

$$\text{Len}(u) = (\text{Len}(P) - 1 - k)$$

Figure 3. The bad-character shift,  $b$  occurs in  $P$ .

Case 2:

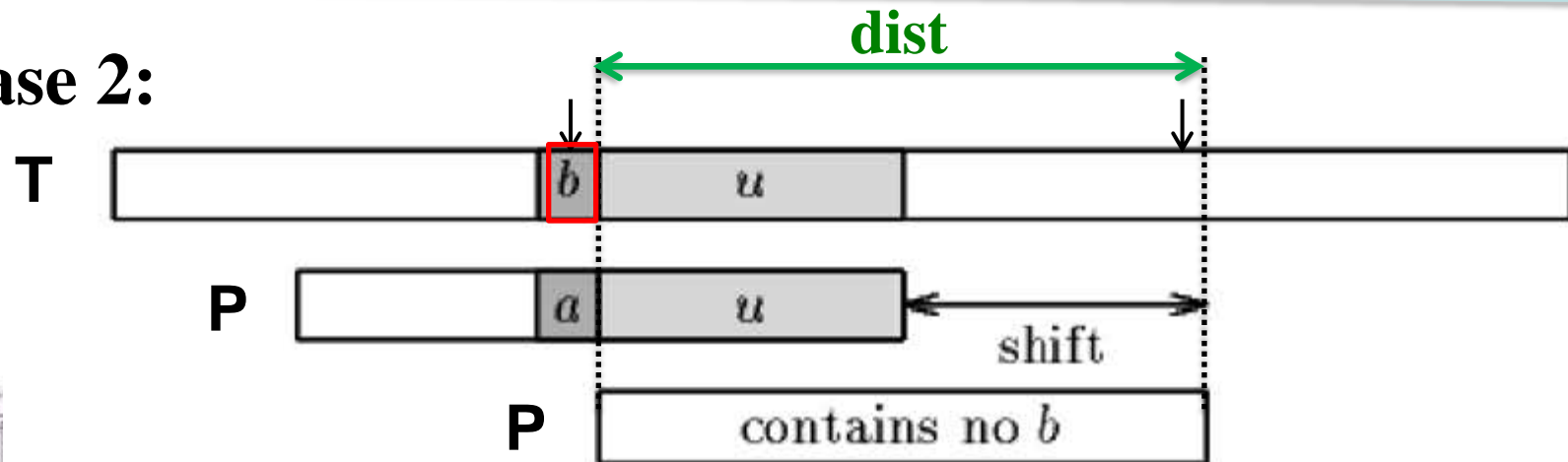
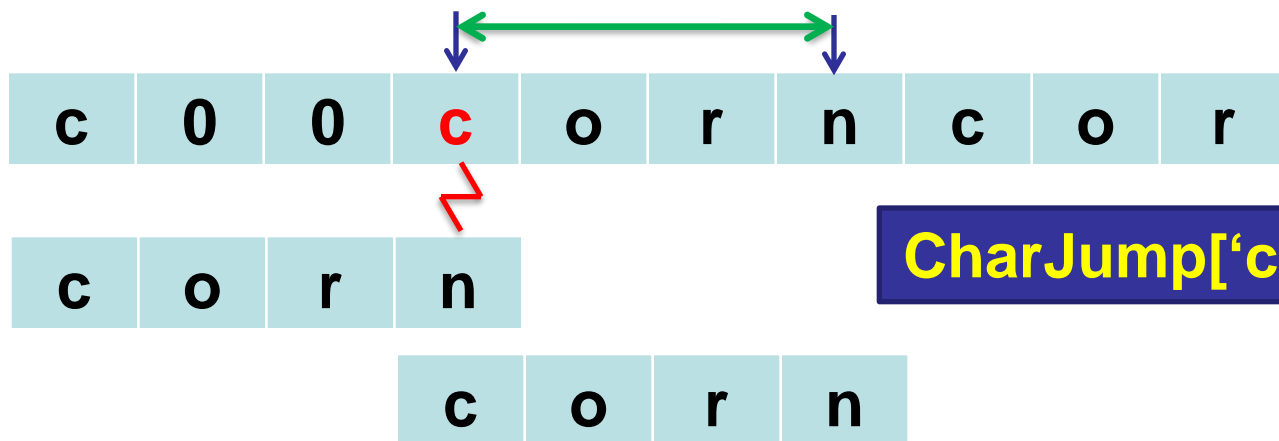


Figure 4. The bad-character shift,  $b$  does not occur in  $P$ .

# CharJump数组

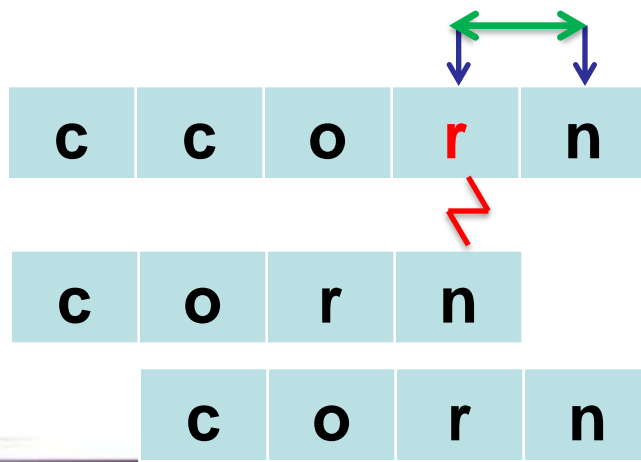
例1：模式串P='corn'

Case1:



CharJump['c']=3

Case2:

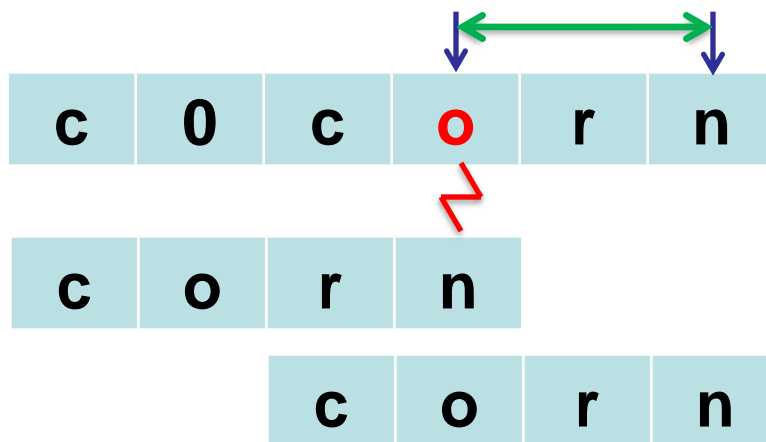


CharJump['r']=1



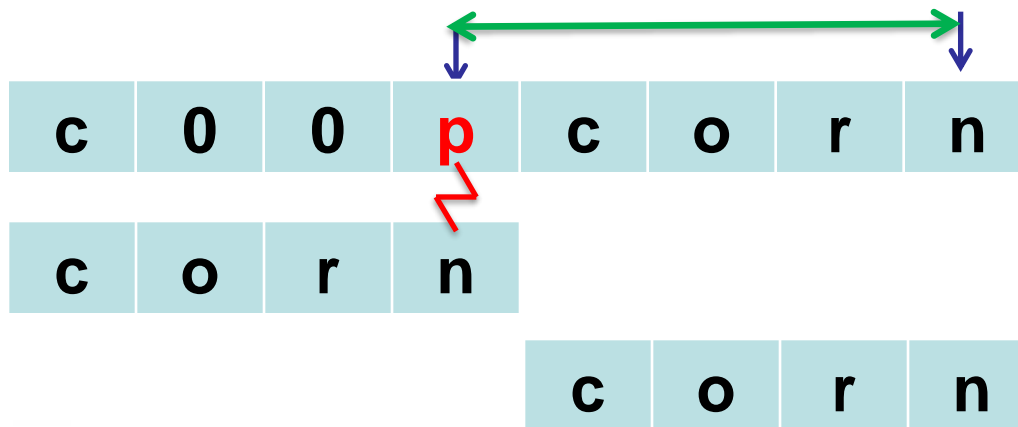
# CharJump数组

Case3:



**CharJump['o']=2**

Case4:

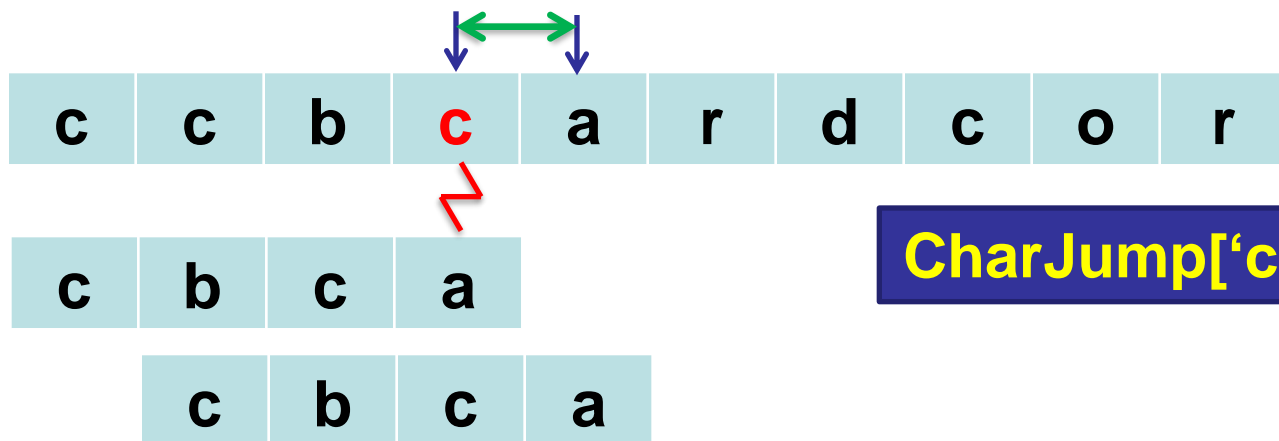


**CharJump['p']=4**

# CharJump数组

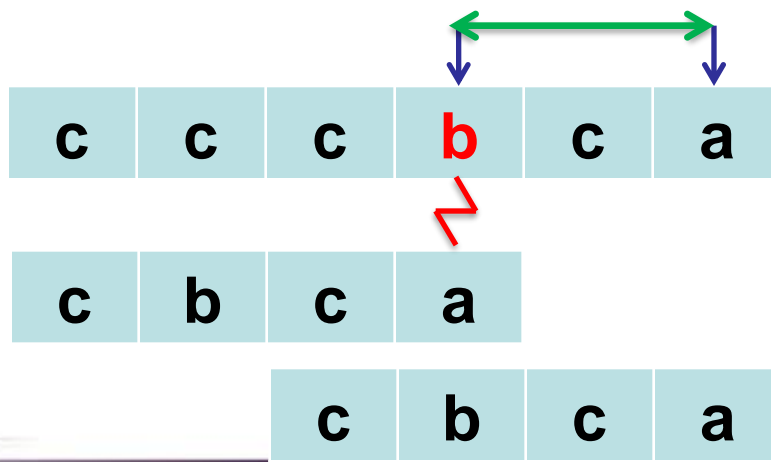
例2：模式串P='cbca'

Case1:



**CharJump['c']=1**

Case2:



**CharJump['b']=2**



# 总结：“坏字符”规则

- “坏字符”规则：假设坏字符为 $x=T[i]$ ，
  - **Case1**：如果坏字符 $x$ 在模式 $P$ 中没有出现，那么从字符 $x$ 开始的 $Len(P)$ 个字符显然不可能与 $P$ 匹配成功，因此，可以使目标串中查找指针直接跳过 $Len(P)$ 个字符。
  - **Case2**：如果坏字符 $x$ 在模式 $P$ 中出现(假设 $P[j]==x$ )，则将目标串中查找指针移动 $CharJump[x]$ ，使得下一轮匹配中将字符 $P[j]$ 与坏字符 $x$ 进行对齐。
- 关键问题：计算 $CharJump[x]$ 
  - 若 $x$ 在 $P$ 中出现，假设 $p[j]==x$ ，则 $CharJump[x]=Len(P)-\max(j)-1$ ；
  - 否则， $CharJump[x]=Len(P)$ ；

例：若 $P$ 为“aedabc”，此时模式串中有重复字符'a'。 $CharJump('a')=2$ ；
- **注意**：模式串 $P$ 右移距离 $shift=CharJump[x] - (Len(P)-1-k)$ 个字符。

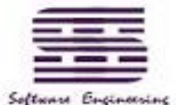
# “坏字符” 数组

```
#define AlphabetSize (UCHAR_MAX + 1) /* For portability */
char *BoyerMoore ( const char * String, /* search for this */
                  const char * Text, /* ...in this text */
                  size_t TextLen ) /* ...up to here. */
{
    /* array of character mismatch offsets */
    unsigned CharJump[AlphabetSize];
    size_t PatLen;
    unsigned u;

    /* Set up and initialize arrays */
    PatLen = strlen ( String );

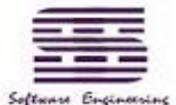
    /* Heuristic #1 -- simple char mismatch jumps ... */
    memset ( CharJump, 0, AlphabetSize * sizeof(unsigned) );
    for ( u = 0 ; u < PatLen; u++ )
        CharJump[((unsigned char) String[u])]
            = PatLen - u - 1;
```

例：若String为"aedabc", CharJump('a')=2;



# 练习：计算CharJump

- 若模式串为"abcd", 则
  - CharJump['a']=3
  - CharJump['b']=2
  - CharJump['c']=1
  - CharJump['d']=0
  - CharJump[all others]=4
- 若模式串为"abcadb", 则
  - CharJump['a']=2
  - CharJump['b']=0
  - CharJump['c']=3
  - CharJump['d']=1
  - CharJump[all others]=6





# 示例：“坏字符”规则的应用

- 例：在“oaks from acorns grow”中查找模式串“corn”的位置。

oaks from acorns grow

Round 1: corn (目标串中指针右移4字符)

Round 2: corn (目标串中指针右移2字符)

Round 3: corn (目标串中指针右移4字符)

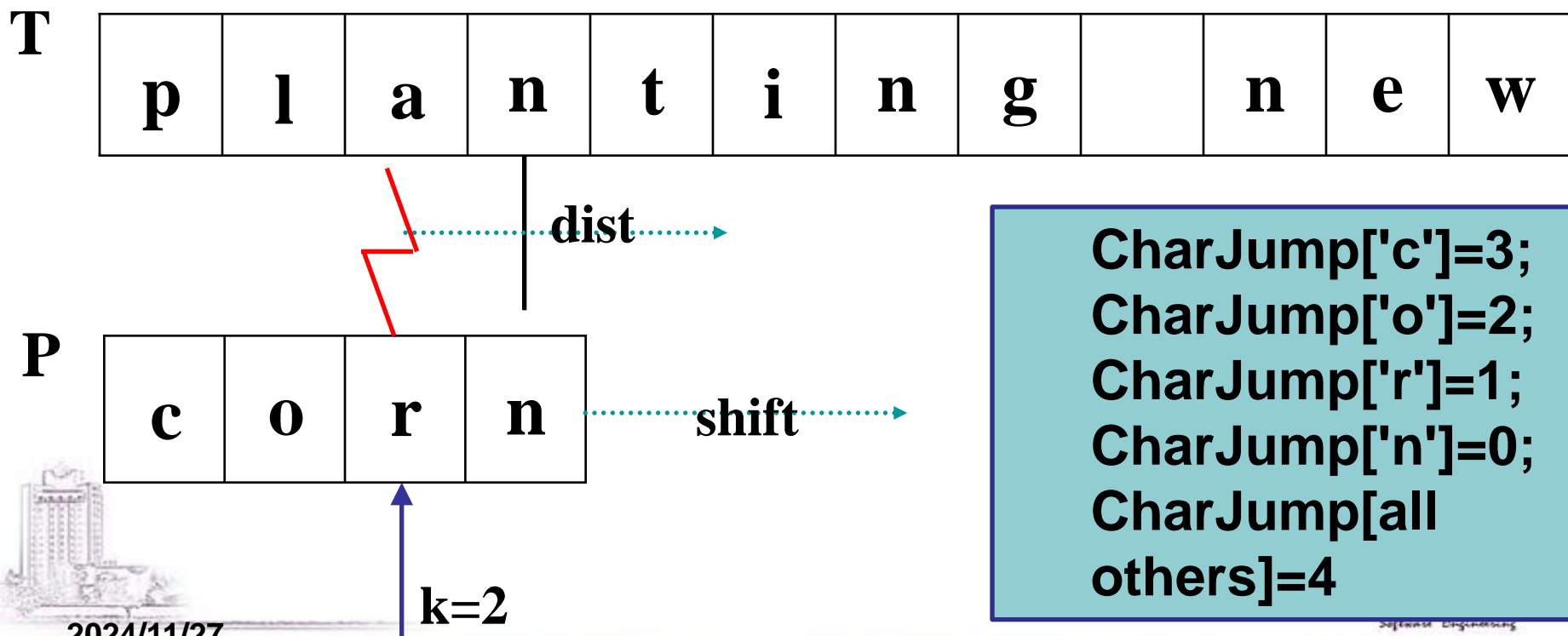
Round 4: corn (目标串中指针右移1字符)

Round 5: corn

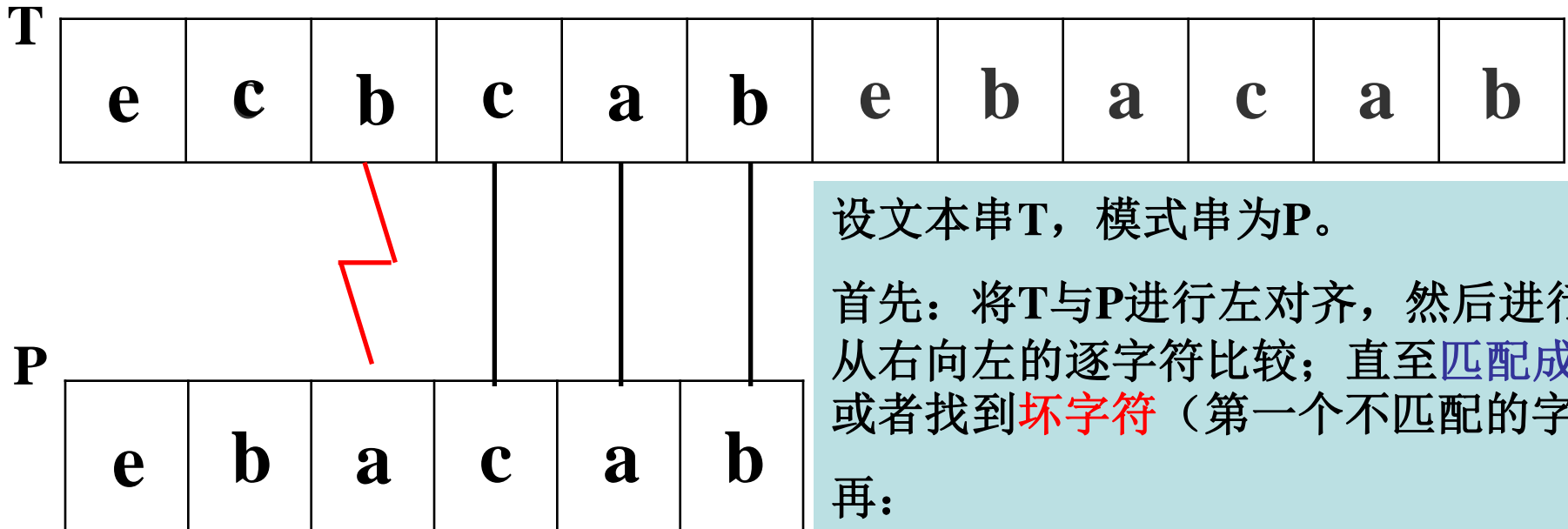
CharJump['c']=3; CharJump['o']=2; CharJump['r']=1;  
CharJump['n']=0; CharJump[all others]=4

# 课堂练习：“坏字符”规则的应用

1、在“planting new”中查找子串“corn”的位置。



2、在“ecbcabebacab”中查找子串“ebacab”的位置。



```
CharJump['e']=5;
CharJump['c']=2;
CharJump['a']=1;
CharJump['b']=0;
CharJump[all others]=6
```

设文本串T，模式串为P。

首先：将T与P进行左对齐，然后进行从右向左的逐字符比较；直至**匹配成功**或者找到**坏字符**（第一个不匹配的字符）

再：

Case1：如果**坏字符**（记为x）在模式P中**没有出现**，那么从字符x开始的Len(P)个字符显然不可能与P匹配成功，因此，可以直接跳过Len(P)个字符以便开始下一轮匹配。

Case2：如果x在模式P中**出现**，则下一轮匹配开始时按照x在模式串中**最右出现的位置**进行对齐。（考虑到x可能在p中多次出现）

2、在“ecbcabebacab”中查找子串“ebacab”的位置。

T

e	c	b	c	a	b	e	b	a	c	a	b
---	---	---	---	---	---	---	---	---	---	---	---



移动距离=0（此时指针移动距离是无效值）

存在好后缀，可以利用“好后缀”规则进行修正，让指针动起来！

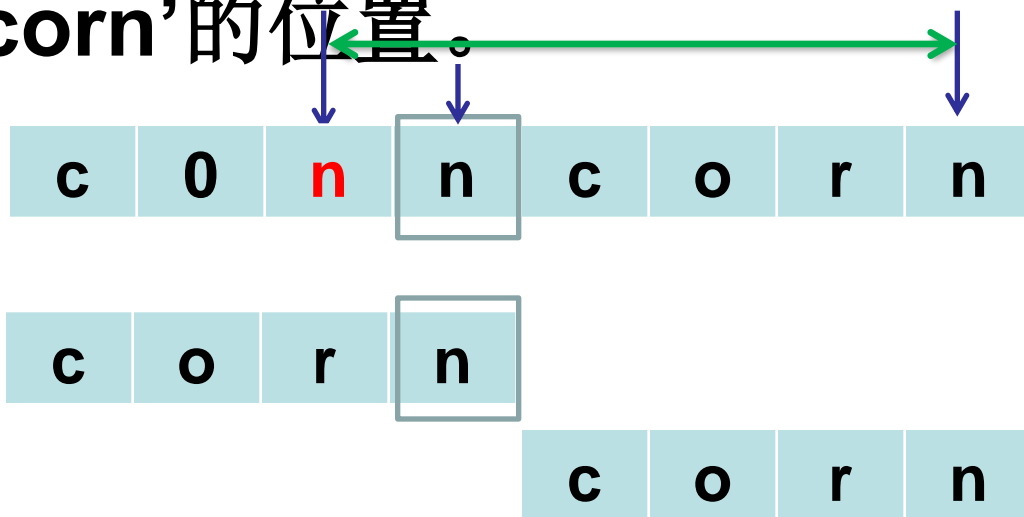
P

e	b	a	c	a	b
---	---	---	---	---	---

**CharJump('b')=0**



### 3、在目标串‘conncorn’中查找模式串 P=‘corn’的位置。



**CharJump['n']=0**

若利用“坏字符规则”，目标串中查找指针的移动距离 $\text{dist}=0$ ，那么此时失配处必定不是在模式串的最后一个字符处，即此时必定存在“好后缀”。

怎么让目标串中查找指针动起来？

利用“好后缀”规则



# “好后缀”规则的三种情况：

1

完全匹配“好后缀”

E D F C A B D A B

0 1 2 3 4 5

C A B D A B

跳 3 位

C A B D A B

跳 4 位

0 1 2 3 4

B C D A B

2

部分匹配“好后缀”

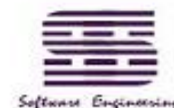
C D A B

跳 4 位

C D A B

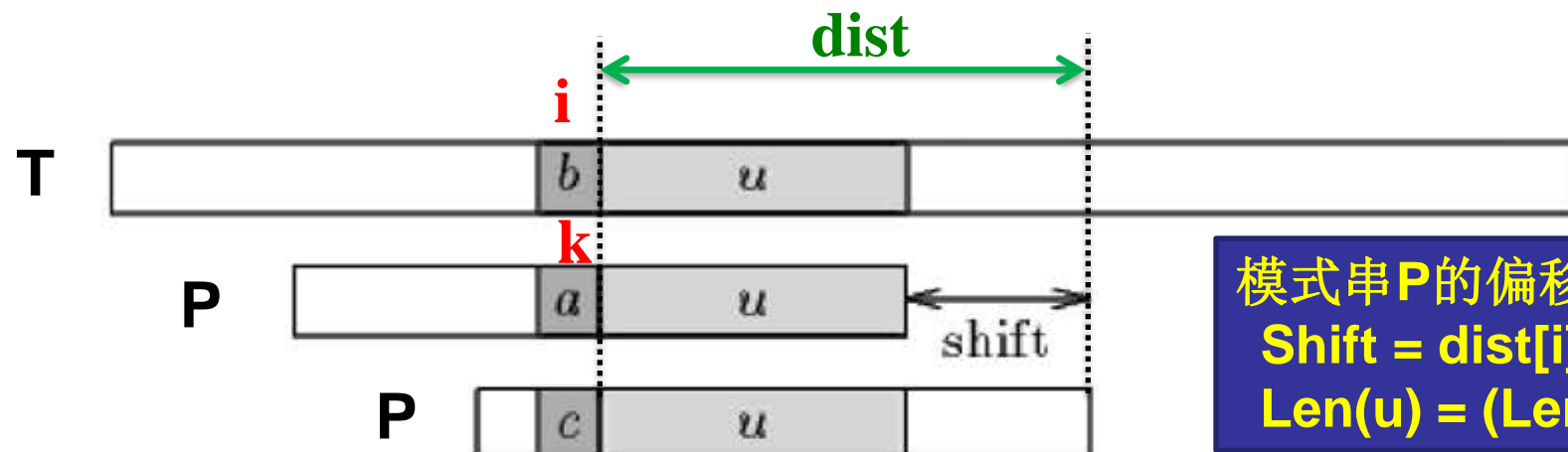
3

未匹配到“好后缀”



# “好后缀”规则的三种情况：

$$\text{dist}[i] = \text{MatchJump}[k]$$



模式串  $P$  的偏移量  
 $\text{Shift} = \text{dist}[i] - \text{Len}(u)$ ,  
 $\text{Len}(u) = (\text{Len}(P) - 1 - k)$

Figure 1. The good-suffix shift,  $u$  re-occurs preceded by a character  $c$  different from  $a$ .

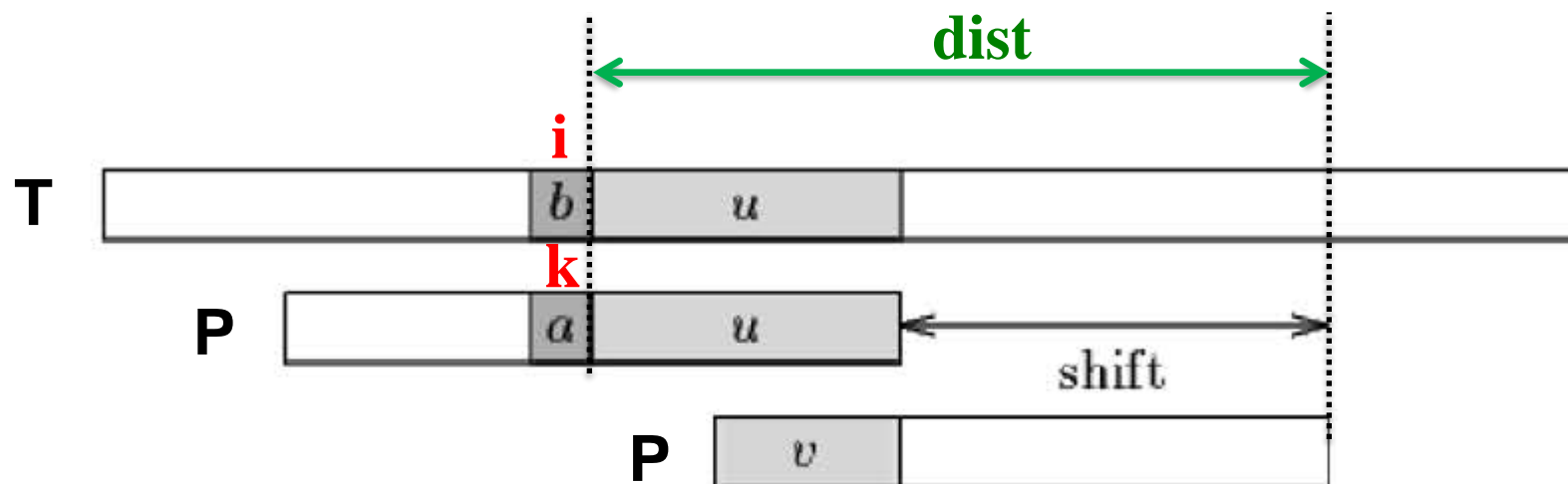


Figure 2. The good-suffix shift, only a suffix of  $u$  re-occurs in  $P$ .



$$\text{dist}[i] = \text{MatchJump}[k]$$

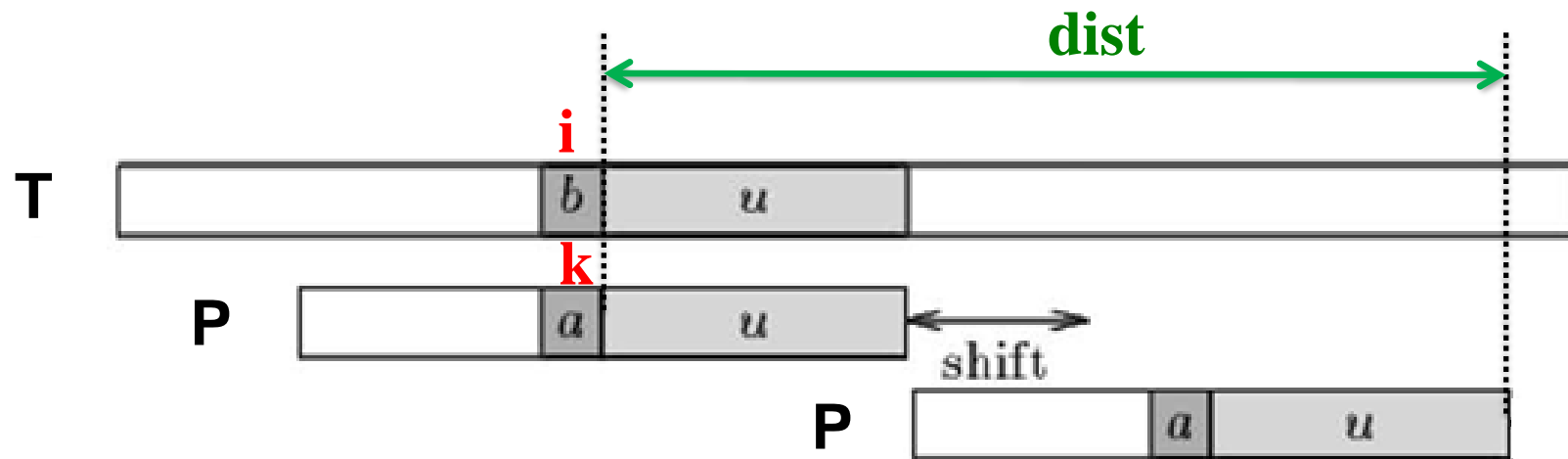


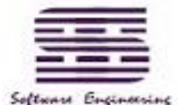
Figure 3. The good-suffix shift, **no re-occurs** in **P**

$$\text{Shift} = \text{dist}[i] - \text{Len}(u),$$

$$\text{Len}(u) = (\text{Len}(P) - 1 - k)$$

Q1: 坏字符可以怎么被用起来?

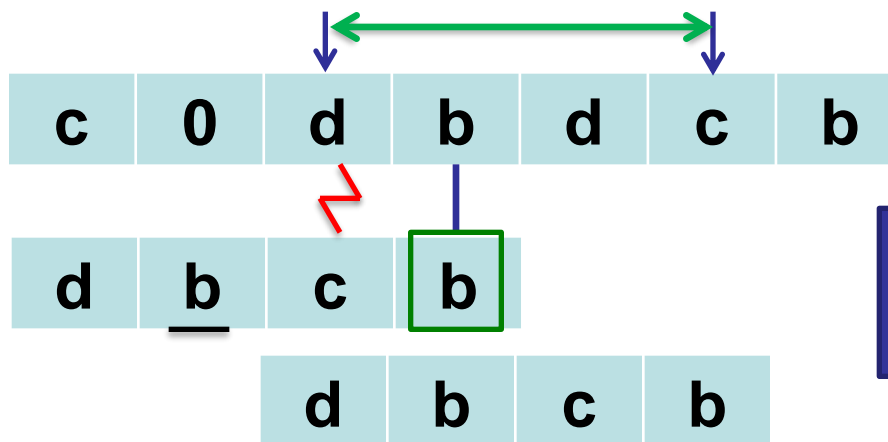
Q2: 若无“好后缀”，该怎么办?





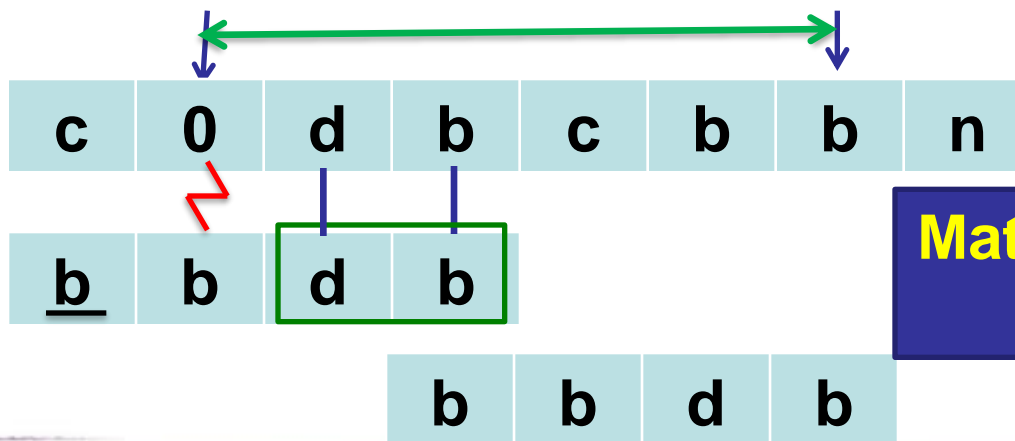
# MatchJump数组

Case1:



**MatchJump[2]=3**  
**=1+2**

Case2:



**MatchJump[2]=5**  
**=2+3**



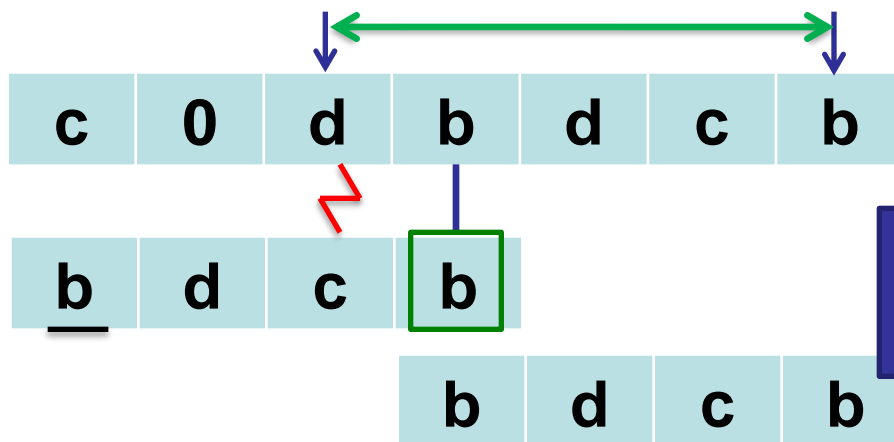
2024/11/27



65  
65

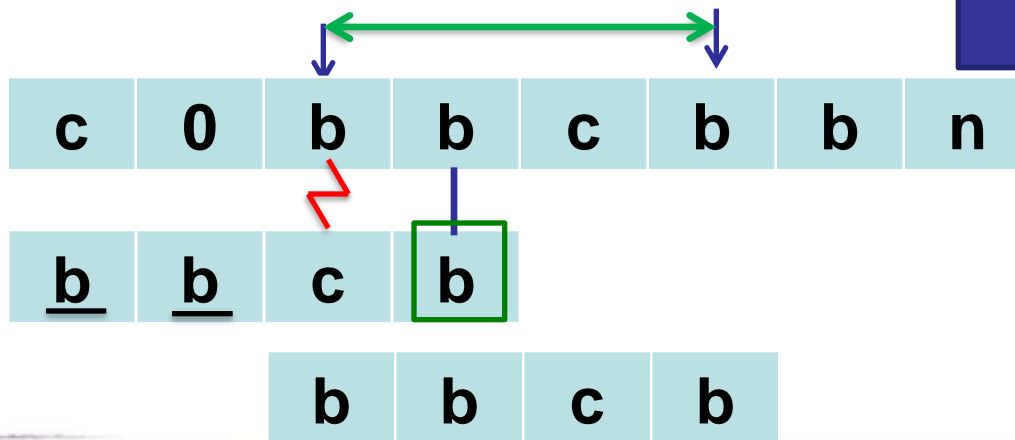
# MatchJump数组

Case3:



**MatchJump[2]=4  
=1+3**

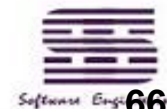
Case4:



**MatchJump[2]=3  
=1+2**



2024/11/27



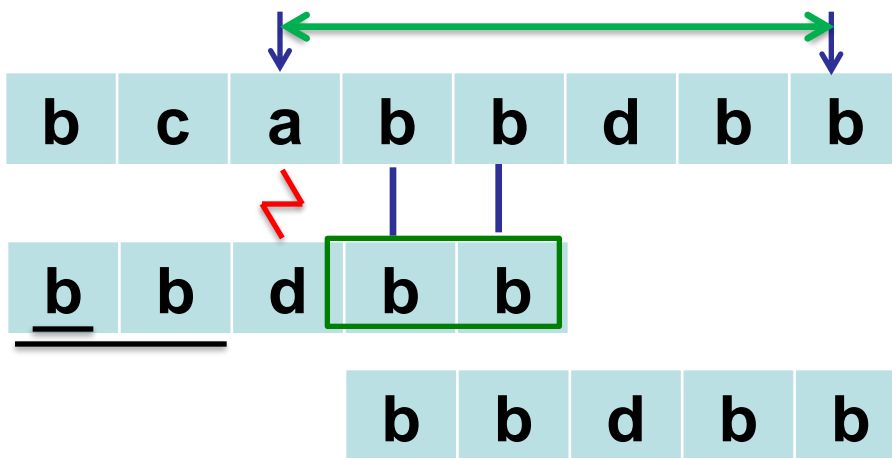
Software Engineering

66

66

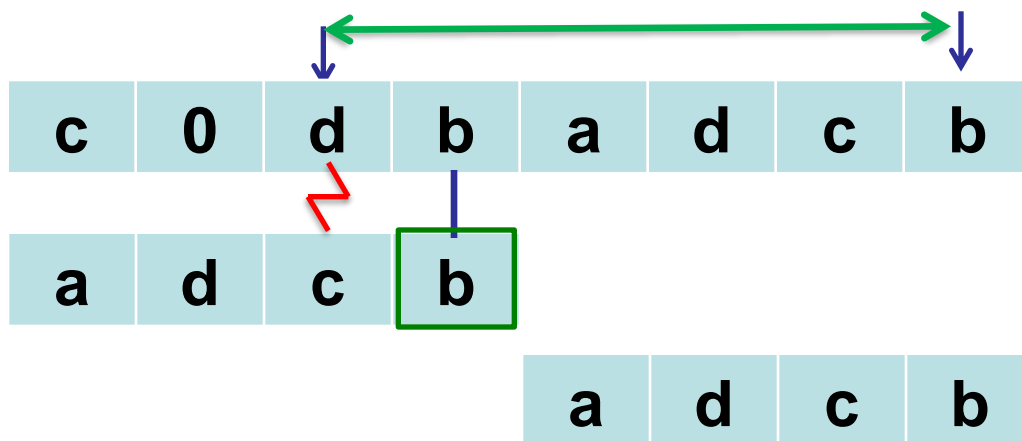
# MatchJump数组

Case5:



**MatchJump[1]=5**  
**=2+3**

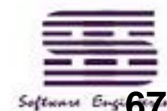
Case6:



**MatchJump[2]=5**  
**=1+4**



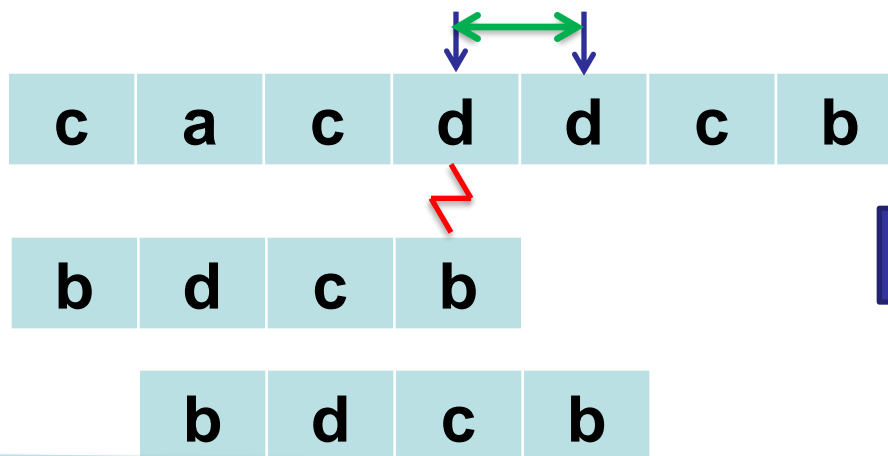
2024/11/27



67  
67

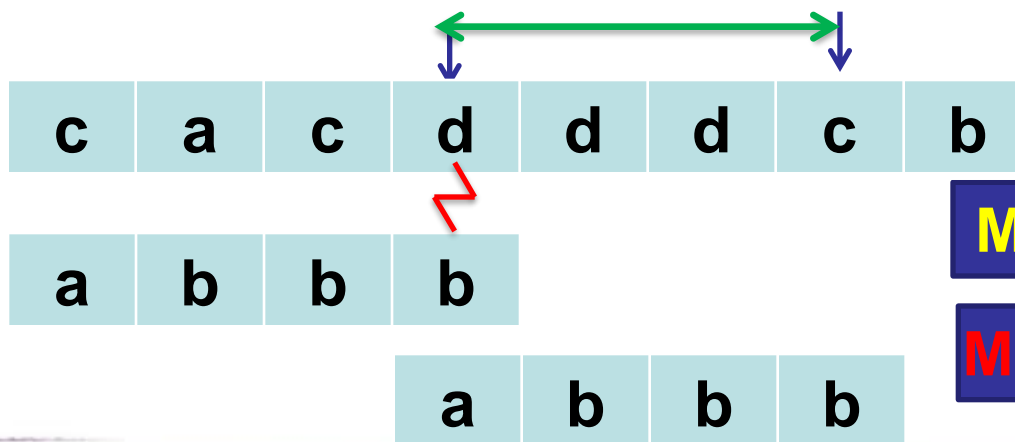
# MatchJump数组

Case7:



**MatchJump[3]=1**

Case8:



**MatchJump[3]=1**

**MatchJump[3]=3?**



2024/11/27

Software Engineering

68

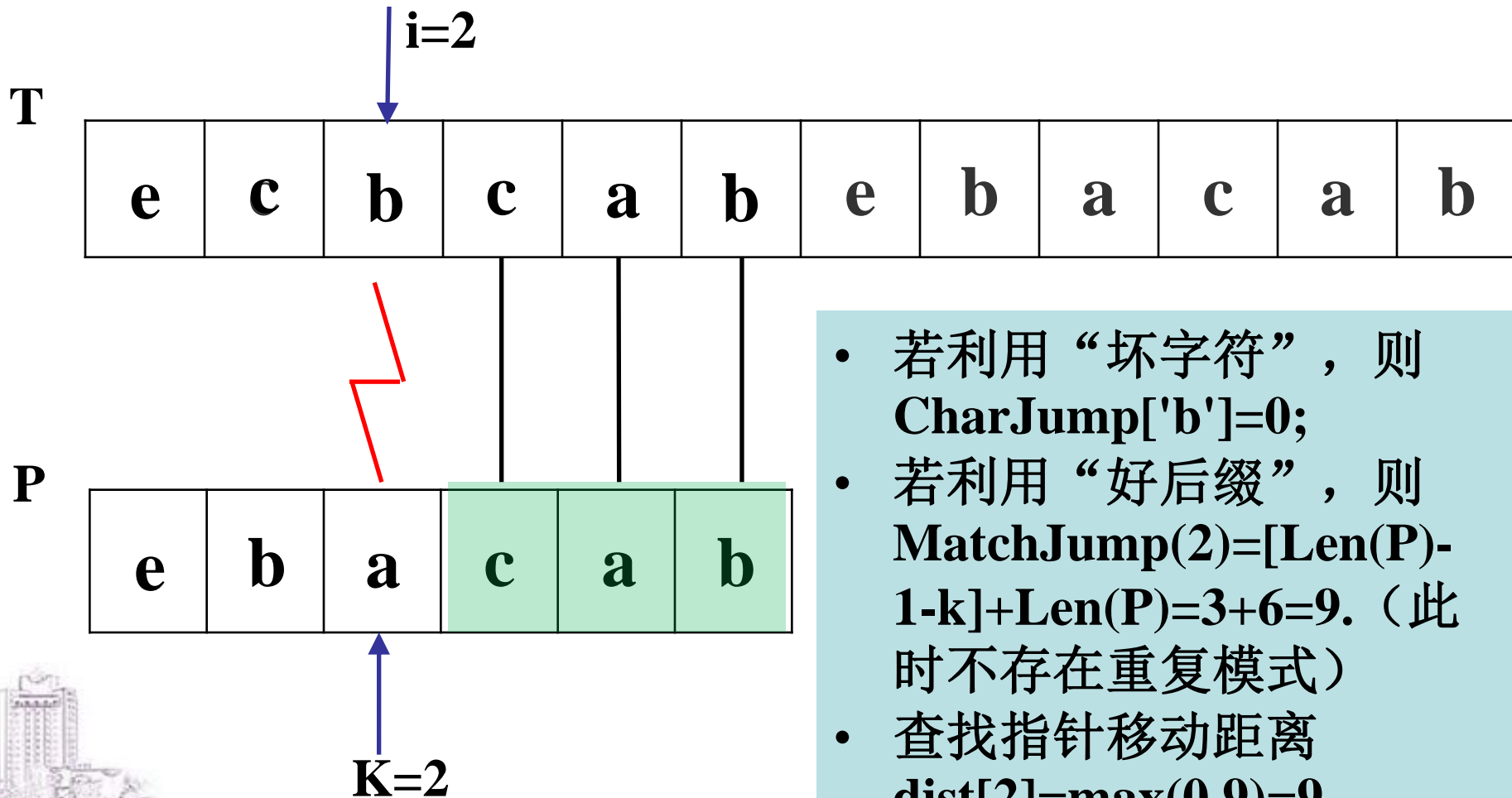
68

# 总结：“好后缀”规则

- 关键问题：如何计算目标串中查找指针的移动距离  
**MatchJump**?
  - T中查找指针的移动距离 $\text{dist}$ =模式串P的移动距离 $\text{shift}$ +好后缀的长度
- **Step1: MatchJump数组初始化:**
  - $\text{MatchJump}[k] = 2 * \text{Len}(p) - k - 1 = \text{Len}(p) - k - 1 + \text{Len}(p);$   
( $k \in [0, \text{Len}(P) - 2]$ )
  - $\text{MatchJump}[\text{Len}(P) - 1] = 1$  (意味着：此时没有好后缀)
- **Step2:**若存在好后缀，则依据“好后缀规则”，重新调整**MatchJump[k]** ( $k \in [0, \text{Len}(P) - 2]$ ) (假定：好后缀P'为  $P[k+1] \dots P[\text{Len}(p) - 1]$ )
  - 情况1：若好后缀P'在P中存在重复模式，且重复模式的前一个字符不等，即 $P[t+1] \dots P[\text{Len}(p) - 1 - k + t] == P[k+1] \dots P[\text{Len}(P) - 1]$ 且 $P[t] \neq P[k]$ ，则 $\text{MatchJump}(k) = [\text{Len}(P) - 1 - k] + \min(k - t)$
  - 情况2：若不满足情况1，且P的前缀为好后缀P'中的某个子后缀P"的重复模式，即 $(P'' = P[t] \dots P[\text{Len}(p) - 1] == P[0] \dots P[\text{Len}(p) - 1 - t])$  ( $t > k + 1$ )，则 $\text{MatchJump}(k) = [\text{Len}(P) - 1 - k] + \min(t)$

(要求计算重复模式的最小间隔，是因为可能有多个重复模式)

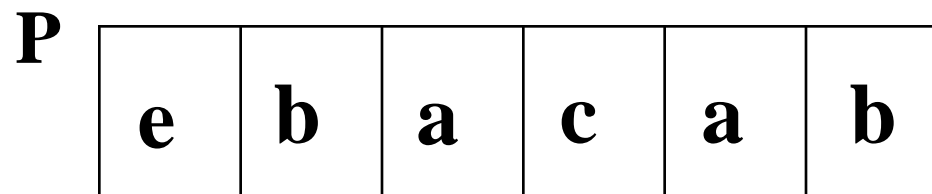
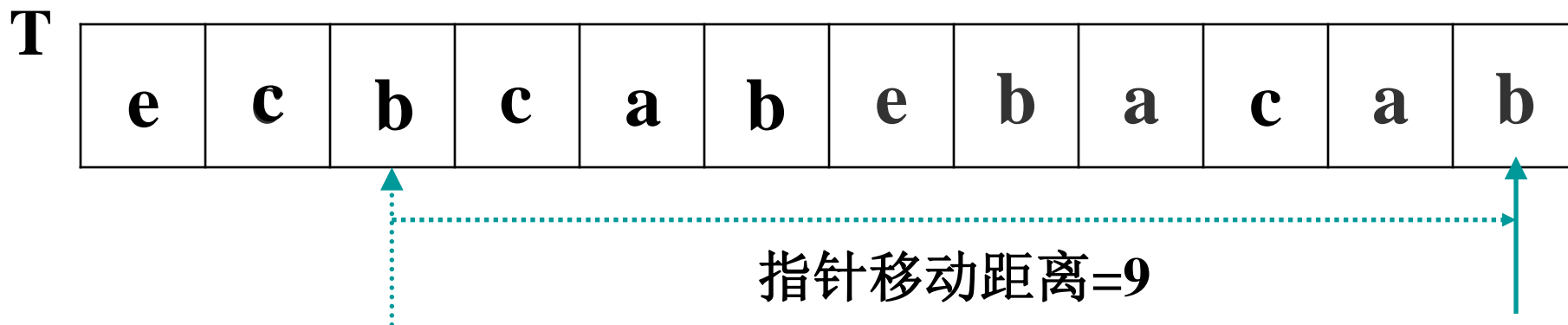
在“ecbcabebacab”中查找子串“ebacab”的位置。



- 若利用“坏字符”，则  $\text{CharJump}['b']=0$ ;
- 若利用“好后缀”，则  $\text{MatchJump}(2)=[\text{Len}(P)-1-k]+\text{Len}(P)=3+6=9$ . (此时不存在重复模式)
- 查找指针移动距离  $\text{dist}[2]=\max(0,9)=9$ .



## Round 2



从右至左，逐个字符匹配，成功！



# 示例2：如何计算MatchJump？

- 若P为：

e	g	f	c	b	b
---	---	---	---	---	---

$(k \in [0, \text{Len}(P) - 1])$

- 分析：模式串p中存在重复模式：后缀子串"b"
- 若好后缀p'中的后缀P''有重复模式，则

- $$\text{MatchJump}(k) = \text{Len}(p) - k - 1 + (t - t')$$

- $\text{MatchJump}(0) = 5 + 6 = 11$

- $\text{MatchJump}(1) = 4 + 6 = 10$

- $\text{MatchJump}(2) = 3 + 6 = 9$

- $\text{MatchJump}(3) = 2 + 6 = 8$

- $\text{MatchJump}(4) = 1 + 1 = 2$  (因为好后缀"b"与"cb"中的"b"存在重复，且  $P[3] \neq p[4]$ )

- $\text{MatchJump}(5) = 1$

好后缀  
的长度





# 示例3： 如何计算MatchJump?

- 若P为:

e	g	f	b	b	b
---	---	---	---	---	---

$k \in [0, \text{Len}(P)-1]$

- 分析：模式串p中存在重复模式：后缀子串"bb"、"b"
- 若好后缀p'中的后缀P''有重复模式，则
- $\text{MatchJump}(k) = \text{Len}(p) - k - 1 + (t - t')$ 
  - $\text{MatchJump}(0) = 5 + 6 = 11$
  - $\text{MatchJump}(1) = 4 + 6 = 10$
  - $\text{MatchJump}(2) = 3 + 6 = 9$
  - $\text{MatchJump}(3) = 2 + 1 = 3$ （因为好后缀"bb"与"fbb"中的"bb"存在重复，且 $P[3] \neq p[2]$ ）
  - $\text{MatchJump}(4) = 1 + 2 = 3$ （因为好后缀"b"与"fb"中的"b"存在重复，且 $P[4] \neq p[2]$ ）
  - $\text{MatchJump}(5) = 1$

好后缀  
的长度

# 示例4：如何计算MatchJump？

- 若P为：

e	b	a	c	a	b
---	---	---	---	---	---



$k \in [0, \text{Len}(P)-1]$

- 分析：模式串p中存在重复模式：后缀子串"b"
- 若好后缀p'中的后缀P''有重复模式，则：

- $\text{MatchJump}(k) = \text{Len}(P) - 1 - k + (t - t')$

- $\text{MatchJump}(0) = 5 + 6 = 11$

- $\text{MatchJump}(1) = 4 + 6 = 10$

- $\text{MatchJump}(2) = 3 + 6 = 9$

- $\text{MatchJump}(3) = 2 + 6 = 8$

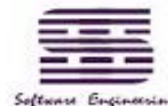
- $\text{MatchJump}(4) = 1 + 4 = 5$ （因为好后缀"b"与"eb"中的"b"存在重复，且 $P[0] \neq p[4]$ ）

- $\text{MatchJump}(5) = 1$

好后缀  
的长度



2024/11/27



Software Engineering

74

# 示例5：如何计算MatchJump？

- 若P为：

e	a	b	b	a	b
---	---	---	---	---	---

$k \in [0, \text{Len}(P)-1]$

好后缀  
的长度

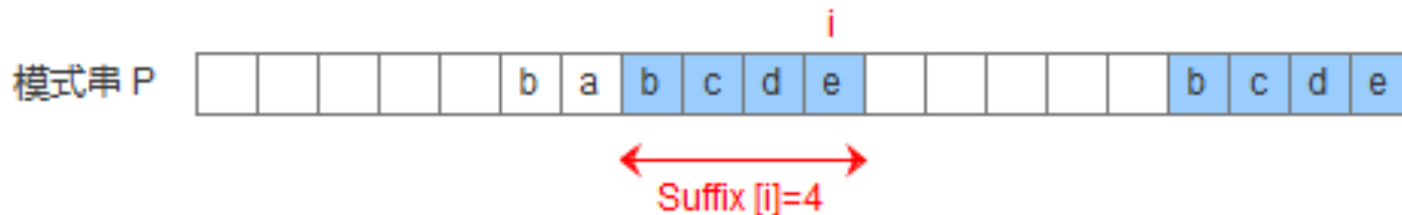
- 分析：模式串p中存在重复模式：后缀子串"ab"、"b"
- 若好后缀p'中的后缀P''有重复模式，则：
- $\text{MatchJump}(k) = \text{Len}(P) - 1 - k + (t - t')$ 
  - $\text{MatchJump}(0) = 5 + 6 = 11$
  - $\text{MatchJump}(1) = 4 + 6 = 10$
  - $\text{MatchJump}(2) = 3 + 6 = 6$
  - $\text{MatchJump}(3) = 2 + 3 = 5$ （因为好后缀"ab"与"eab"中的"ab"存在重复，且 $P[0] \neq p[3]$ ）
  - $\text{MatchJump}(4) = 1 + 2 = 3$ （因为好后缀"b"与"bb"中的"b"存在重复，且 $P[2] \neq p[4]$ ）
  - $\text{MatchJump}(5) = 1$

# 示例6：如何计算MatchJump?

- 若P为"acebabaceb":
  - 分析：模式串p中存在后缀子串"aceb"、"ceb"、"eb"、"b"的重复模式
  - 则：
    - $\text{MatchJump}(0)=9+6=15$  (P的前缀"aceb"与好后缀中"cebabaceb"的子后缀"aceb"存在重复。好后缀长度为9，重复模式间隔为6)
    - $\text{MatchJump}(1)=8+6=14$
    - $\text{MatchJump}(2)=7+6=13$
    - $\text{MatchJump}(3)=6+6=12$
    - $\text{MatchJump}(4)=5+6=11$
    - $\text{MatchJump}(5)=4+6=10$
    - $\text{MatchJump}(6)=3+10=13$  (k=6,意味着在“acebabaceb”带下划线的字符处匹配失败，观察“acebabaceb”，并不满足情况1，(将黄色字符对齐是多余的)按照情况2，找好后缀的子串，仍然没有，故取其初值3+10)
    - $\text{MatchJump}(7)=2+10=12$  (与MatchJump(6)的计算原理类似)
    - $\text{MatchJump}(8)=1+4=5$  (好后缀为"b"，P中存在1个好后缀的重复模式，如“acebabaceb”，且P[8] ≠ P[4]，重复模式间隔为4)
    - $\text{MatchJump}(9)=1$

# “好后缀”数组 (shift距离)

为了实现好后缀规则，需要定义一个数组suffix[]，其中  
**suffix[i] = s** 表示以i为边界，与模式串后缀匹配的最大长度，  
即：满足 $P[i-s, i] == P[m-s, m]$ 的最大长度s。



```
void suffixes(char *x, int m, int *suff){
    suff[m-1]=m;
    for (i=m-2; i>=0; --i){
        q=i;
        while(q>=0&& x[q]==x[m-1-i+q]) --q;
        suff[i]=i-q;  } }
```



# “好后缀”数组 (shift距离)

初始化情况（好后缀不重复出现）

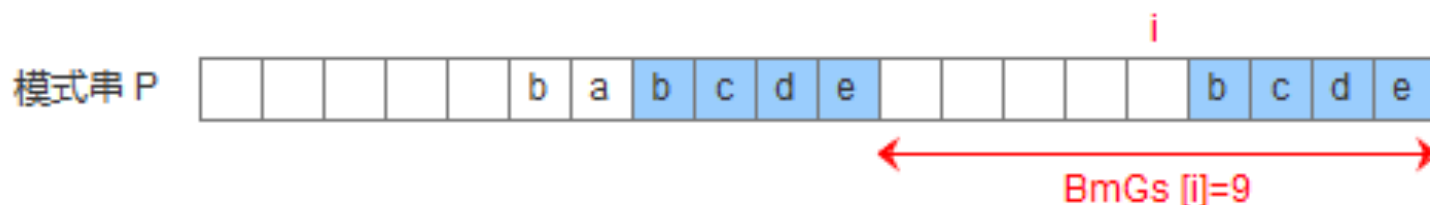


初始化，BmGs[i]均为P串的大小(Len(P))

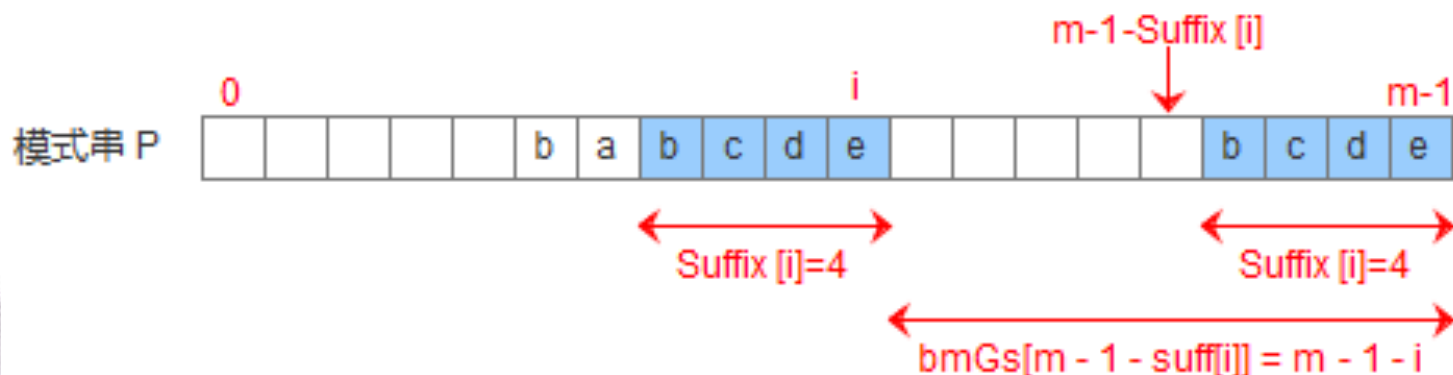


# “好后缀”数组 (shift距离)

有重复后缀 (情况1)



计算视角? : 当P串指针为 $i$ 时, mismatch



## “好后缀” 数组 (shift距离)

## 有最大后缀子串（情况2）



对于上述情况， $BmGs[2\sim17] = 18$  也就是空白格中的值都应该是18





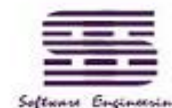
# “好后缀”数组 (shift距离)

有最大后缀子串（情况2），如何计算



字符串匹配： $x[i-(suff[i]-1) \sim i] = x[m-1-suff[i] \sim m-1]$ ，  
而前缀必须从0开始：也就是  $x[i-(suff[i]-1) \sim i] = x[0, i]$   
故  $suff[i] = i+1$  时，满足第二种情况。

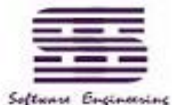
找到满足  $suff[i] = i+1$  的最大  $i$  值



# “好后缀”数组 (shift距离)

```
void preBmGs(char *x, int m, int bmGs[]) {  
    int i, j, suff[XSIZE];  
    suffixes(x, m, suff);  
    for (i = 0; i < m; ++i) bmGs[i] = m; // 初始化情况, 无匹配  
    j = 0;  
    for (i = m - 1; i >= 0; --i)  
        if (suff[i] == i + 1)  
            for (; j < m - 1 - i; ++j)  
                if (bmGs[j] == m) // 值为初值时更改  
                    bmGs[j] = m - 1 - i; // 子缀匹配  
  
    for (i = 0; i <= m - 2; ++i)  
        bmGs[m - 1 - suff[i]] = m - 1 - i; // 好后缀匹配  
}
```

<https://www.cnblogs.com/xubenben/p/3359364.html>



# BM算法的应用

- 目标串T: **GCATCGCAGAGAGTATACAGTACG**
- 模式串P: **GCAGAGAG**

x	A	C	G	其他字符
CharJump(x)	1	6	0	8

- 分析: 模式串p中存在后缀子串“G”,  
“AG”, “GAG”, “AGAG”的重复模式

K	0	1	2	3	4	5	6	7
P[k]	G	C	A	G	A	G	A	G
MatchJump[k]	7+7	6+7	5+7	4+2	3+7	2+4	1+7	1

### First attempt

G C A T C G C A G A G A G T A T A C A G T A C G

1

G C A G A G A G

Shift by: 1 ( $bmGs[7] = bmBc[A] - 8 + 8$ )

### Second attempt

G C A T C G C A G A G A G T A T A C A G T A C G

3 2 1

G C A G A G A G

Shift by: 4 ( $bmGs[5] = bmBc[C] - 8 + 6$ )



Third attempt

G C A T C G C A G A G A G T A T A C A G T A C G

8 7 6 5 4 3 2 1

G C A G A G A G



例题：给定一个非空的字符串  $s$ ，检查是否可以通过由它的一个子串重复多次构成。

示例 1:

输入:  $s = \text{"abab"}$

输出: `true`

解释: 可由子串 `"ab"` 重复两次构成。

示例 2:

输入:  $s = \text{"aba"}$

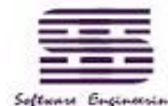
输出: `false`

示例 3:

输入:  $s = \text{"abcabcabcabc"}$

输出: `true`

解释: 可由子串 `"abc"` 重复四次构成。(或子串 `"abcabc"` 重复两次构成。)



# BM算法的性能点评

- **BM算法的准备时间：与Patlen+AlphabetSize成正比。**
- **BM算法利用了2个跳转表：CharJump和MatchJump，跳过了目标字符串中某些字符，不需要对目标字符串中所有字符进行逐一比较。**
- **通常模式串越长，BM算法速度越快。因为每一次失败的匹配尝试，BM算法都能够使用这些信息（好后缀和坏字符）来排除尽可能多的无法匹配的位置。**

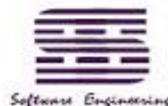


# 附：作业7（课后）？

- 利用**BM**算法和**KMP**算法完成给定的字符串查找。
- 总结这两类字符串查找算法的特点、并进行性能比较。



2024/11/27



88