



Beijing Normal University
School of Artificial Intelligence

第4章 字符串

郑新、徐鹏飞、李健
北京师范大学人工智能学院
2025-2026学年 第一学期

考核要点

■ 考核大纲

- 字符串的模式匹配

■ 复习要点

- 理解字符串的基本概念 (A)
- 熟练掌握字符串的顺序实现和链接实现 (A)
- 运用字符串求解经典问题 (B)
 - 了解朴素模式匹配 (BF暴力匹配) 算法
 - 重点掌握KMP匹配算法的原理

- 4.1 问题引入：模式匹配
- 4.2 字符串的定义与结构
- 4.3 字符串的存储实现
- 4.4 字符串的模式匹配
- 4.5 拓展延伸
- 4.6 应用场景：基因测序

4.1 问题引入：模式匹配

- **查找功能：**文本编辑工具中，给定一段文本，用户提供特定的关键词，找出这个关键词在文本中出现的位置，就是字符串匹配问题。



- 用户给定的关键词或字符串被称为**模式串**，段落文本成为**目标串**，因此字符串匹配又称为**模式匹配**。

- 4.1 问题引入：模式匹配
- 4.2 字符串的定义与结构
- 4.3 字符串的存储实现
- 4.4 字符串的模式匹配
- 4.5 拓展延伸
- 4.6 应用场景：基因测序

4.2 字符串的定义与结构

- 字符串（String），简称串，由零个或多个字符组成的有限序列，记为 $s = "s_1s_2...s_n"$ ($n > 0$) 或 \emptyset
 - 一种在数据元素的组成上具有一定约束条件的线性表，即要求组成线性表的所有数据元素都是字符
 - 其中s为字符串名
 - 双引号为字符串的定界符，双引号之间的内容是字符串的值
 - s_i ($1 \leq i \leq n$) 可以是字母、数字或其他字符
 - n为字符串长度。 \emptyset 表示空串，空串的长度为0
 - 一个或多个空格组成串称为空格串，空格串不是空串
 - 例如 " "和""分别表示长度为1的空格串和长度为0的空串。
 - 两个串长度相等，且每个对应位置的字符都相等时，两个串相等

4.2 字符串的定义与结构

■ 空串、空格串、空白串的区别

- 空串是“无字符”，长度为 0；
- 空格串是“仅含空格”，长度 ≥ 1 ，属于空白串的子集；
- 空白串是“含任意空白字符”，长度 ≥ 1 ，范围包含空格串。
 - 空格 (' ')、制表符 ('\t')、换行符 ('\n')、回车符 ('\r') 等所有空白字符

对比维度	空串 (Empty String)	空格串 (Space String)	空白串 (Whitespace String)
定义	不含任何字符的字符串	仅由空格字符 (' ', ASCII 32) 组成的字符串	由任意一种或多种空白字符组成的字符串
组成字符	无任何字符	仅包含空格字符 (' ')	包含空格 (' ')、制表符 ('\t')、换行符 ('\n')、回车符 ('\r') 等所有空白字符
长度	长度为 0 ($\text{len}("") = 0$)	长度 ≥ 1 (至少包含 1 个空格)	长度 ≥ 1 (至少包含 1 个空白字符)
示例	"" (双引号中无任何内容)	" " (3 个连续空格)	" \t\n" (1 个空格 + 1 个 Tab + 1 个换行)
本质特征	表示“无字符”的状态	属于“空白串”的特殊子集 (仅含空格)	包含所有空白字符，范围最广

4.2 字符串的定义与结构

- **子串 (substring)**：串中任意个连续字符组成的子序列称为该串的子串，包含子串的串相应地称为主串。
 - 子串在主串中的位置等于子串第一个字符在主串中的位置
 - 例如，设A和B分别为
 $A = \text{"This is a string"}$ $B = \text{"is"}$
则 B 是 A 的子串，A 为主串。B 在 A 中出现了两次，首次出现所对应的主串位置是2（从0开始）。因此，称 B 在 A 中的位置为2。
 - 如长度为4的字符串 $s = \text{"abcd"}$ ，其子串包括 "a"、"ab"、"abc"、"abcd"、"b"、"bc"、"bcd"、"c"、"cd"、"d"、""（空字符串）。
 - 特别地，空串是任意串的子串，任意串是其自身的子串。

4.2 字符串的定义与结构

- 通常在程序中使用的串可分为两种：串变量和串常量。
 - 串常量在程序中只能被引用但不能改变它的值，即只能读不能写。通常串常量是由直接量来表示的，例如语句 `Error (“overflow”)` 中“overflow”是直接量。但有的语言允许对串常量命名，以使程序易读、易写。如C中可定义
`char path[] = “dir/bin/appl”;`
 - 这里path存储的是一个串常量。串变量和其它类型的变量一样，其取值可以改变。

4.2 字符串的定义与结构

■ 字符串的抽象数据类型定义

ADT String {

数据对象: $D = \{ s_i \mid s_i \in \text{CharacterSet}, i=1, 2, \dots, n, n>0 \}$ 或 \emptyset , 表示空字符串

数据关系: $R = \{ \langle s_{i-1}, s_i \rangle \mid s_{i-1}, s_i \in D, i=2, \dots, n, n>0 \}$

基本操作:

InitStr(s): 初始化一个空的字符串s, 字符串最大长度为kMaxSize。

StrCopy(s): 返回复制字符串s得到的字符串。

StrIsEmpty(s): 判断字符串s是否是空串。返回一个布尔值, 若字符串s是空串, 返回true; 否则返回false。

StrInsert(s, pos, t): 在字符串s的pos位置处插入字符串t, 并返回插入后的字符串s。

StrRemove(s, pos, len): 删除字符串s中从pos位置开始的长度为len的子串, 返回删除后的字符串s。

SubString(s, pos, len): 返回字符串s从pos位置开始长度为len的子串。

StrLength(s): 返回字符串s的长度。

StrConcat(s,t): 返回字符串s和t连接而成的新串s。

StrCompare(s,t): 返回字符串s和t的大小关系。若s>t, 返回+1; 若s=t, 返回0; 若s<t, 返回-1。

PatternMatch(s,t): 返回字符串s中首次出现字符串t的位置, 若字符串s没有出现字符串t, 则返回NIL。

Replace(s, sub_s, t): 将字符串s中的所有子串sub_s替换为字符串t。

}

- 4.1 问题引入：模式匹配
- 4.2 字符串的定义与结构
- 4.3 字符串的存储实现
- 4.4 字符串的模式匹配
- 4.5 拓展延伸
- 4.6 应用场景：基因测序

4.3 字符串的存储实现

字符串存储表示

顺序存储结构

将所有数据元素存放在一段连续的存储空间中，数据元素的存储位置反应了它们之间的逻辑关系

链接存储结构

逻辑上相邻的数据元素不需要在物理位置上也相邻，数据元素的存储位置可以是任意的

4.3.1 字符串的顺序存储实现

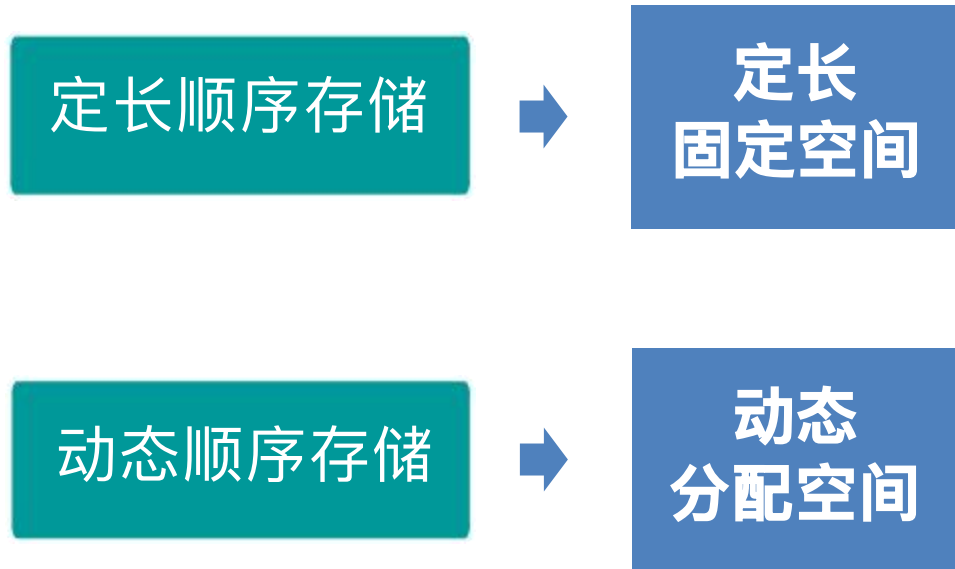
- 字符串的顺序存储结构是用一组地址连续的存储单元来存储串中的字符序列，也是最为常见的字符串存储结构。

- 例：字符串"abcdef"的顺序存储形式

0	1	2	3	4	5	6
a	b	c	d	e	f	\0

- 注意：一般在程序设计语言中，字符串会有结束符，如C语言和C++语言中字符串的结束符为'\0'，结束符不计入字符串长度，但要占存储空间。
- 根据是否预先确定串的存储空间大小，可将字符串的顺序存储结构分为**定长顺序存储**和**动态顺序存储**两种类型。

4.3.1 字符串的顺序存储实现



```
#define MaxLen 255 // 1. 定义字符串的最大长度
typedef struct SString { // 2. 定义顺序串
    char ch[MaxLen]; // 2.1 定义字符串序列
    int length; // 2.2 定义串内字符长度变量
}SString;
```

```
#typedef struct HString { // 1. 定义动态存储的堆字符串
    char *ch; // 2. 按串长分配存储区
    int length; // 3. 定义串内字符长度变量
}HString;
```

存储方式	空间分配时机	位置	空间灵活性	内存管理	适用场景
定长顺序存储	编译时固定分配	栈	大小固定（由数组定义决定），超出则栈溢出	函数 / 作用域结束时自动回收	长度固定且较小的字符串
动态顺序存储	运行时动态分配	堆	可按需调整大小（通过realloc 等）	需手动分配（malloc/new）和释放（free/delete）	长度不确定或较大的字符串

4.3.1 字符串的顺序存储实现

■ 字符串插入

在字符串s的pos位置插入字符串t

- 将字符串s从pos位置开始整体后移，为字符串t插入提供空间
- 将字符串t中的字符逐一插入到字符串s从pos开始的空间中
- 更新字符串s的length属性

■ 时间复杂度

最坏情况是将t插在s的头，时间复杂度是 $O(n+m)$ ，其中n和m分别是两个字符串的长度。

算法4-1 字符串插入操作StrInsert(s, pos, t)

输入：字符串s，要插入的位置 $pos \geq 1$ ，需要插入的字符串t

输出：完成插入后的字符串s。若插入后的字符串长度大于kMaxSize，则直接退出

```
n ← s.length
m ← t.length
if n+m ≤ kMaxSize then
| for i ← n-1 downto pos-1 do
| | s.data[i+m] ← s.data[i] //将数组下标pos-1开始的子串后移，给t留出空位
| end
| for i ← 0 to m-1 do
| | s.data[pos-1+i] ← t.data[i] //将t插入s
| end
| s.length ← n + m //更新s的长度
else
| 长度超限，退出
end
```

4.3.1 字符串的顺序存储实现

```
void StrInsert(String s, Position pos, String t)
{ /* pos从1开始*/
    int n, m;
    Position i;
    n = s->length;
    m = t->length;
    if ((n+m)<kMaxSize) { /* 如果插入后不超长 */
        for (i=n; i>=(pos-1); i--) { /* i=n时将结束符后移 */
            s->data[i+m] = s->data[i]; /*将下标pos-1开始的子串后移, 给t留出空位*/
        }
        for (i=0; i<m; i++) {
            s->data[pos-1+i] = t->data[i]; /* 将t插入s */
        }
        s->length = n + m;
    }
    else {
        printf("错误: 插入将导致字符串长度超限。\\n");
    }
}
```


4.3.1 字符串的顺序存储实现

■ 字符串删除

将字符串s从pos位置开始删除长度为len的子串

- 将字符串s中pos+len后面的部分逐位向前移动
- 更新字符串s的length属性

■ 时间复杂度

最坏情况是将s开头的len个字符删除，时间复杂度是 $O(n)$ ，其中n是s的长度。

算法4-2 字符串删除操作StrRemove(s, pos, len)

输入：字符串s，要删除的位置 $\text{pos} \geq 1$ ，删除的字符个数len

输出：删除子串后的字符串s

```
n ← s.length
if pos+len-1 < n then
| for i ← pos+len-1 to n-1 do
| | s.data[i-len] ← s.data[i]
| end
| s.length ← n - len
else //从数组下标pos-1开始的所有字符都删掉
| s.length ← pos-1
end
```

4.3.1 字符串的顺序存储实现

```
void StrRemove(String s, Position pos, int len)
{ /* pos从1开始 */
    int n;
    Position i;
    n = s->length;
    if ((pos+len-1)<=n) {
        for (i=pos+len-1; i<=n; i++) { /* i=n处理结束符 */
            s->data[i-len] = s->data[i];
        }
        s->length = n - len;
    }
    else {
        s->data[pos-1] = '\0'; /* 从数组下标pos-1开始的所有字符都删掉 */
        s->length = pos-1;
    }
}
```

4.3.1 字符串的顺序存储实现

■ 字符串截取

返回字符串s从pos位置开始长度为len的子串

- 构造新串，将字符串s从pos之后的len个字符一一赋值给新串

■ 时间复杂度

该操作仅与子串长度有关，为 $O(len)$ 。

算法4-3 字符串截取子串操作SubString(s, pos, len)

输入：字符串s，开始截取的位置 $pos \geq 1$ ，截取的字符个数len

输出：截取的子串

InitStr(sub_s) //初始化新串sub_s

$n \leftarrow s.length$

$i \leftarrow 0$

while $pos-1+i < n$ 且 $i < len$ do //若s从pos-1开始不到len个字符，就截取到s的末尾为止

 | $sub_s.data[i] \leftarrow s.data[pos-1+i]$ //将s串从pos-1之后的len个字符复制到sub_s

 | $sub_s.length \leftarrow sub_s.length + 1$

 | $i \leftarrow i+1$

end

return sub_s

4.3.1 字符串的顺序存储实现

```
String SubString(String s, int pos, int len)
{ /* pos从1开始 */
    int n;
    Position i;
    String sub_s;
    sub_s = InitStr(); /* 初始化新串sub_s */
    n = s->length;
    i = 0;
    while (((pos-1+i)<=n) && (i<len)) {
        /* 若s从pos-1开始不到len个字符, 就截取到s的末尾为止 */
        /* 将s串从pos-1之后的len个字符复制到sub_s */
        sub_s->data[i] = s->data[pos-1+i];
        sub_s->length++;
        i++;
    }
    sub_s->data[i] = '\0';
    return sub_s;
}
```

4.3.1 字符串的顺序存储实现

■ 字符串连接

将字符串t连接在字符串s末尾

- 将字符串t中的字符一一赋值到s末尾
- 更新s的长度

■ 时间复杂度

该操作仅与t的长度有关，为 $O(m)$ 。

算法4-4 字符串连接操作StrConcat(s,t)

输入：字符串s，字符串t

输出：返回字符串s后面联接t而成的新串。若结果长度大于kMaxSize, 则退出。

```
n ← s.length  
m ← t.length  
if n+m ≤ kMaxSize then  
| for i ← 0 to m-1 do  
| | s.data[n+i] ← t.data[i]  
| end  
| s.length ← n+m  
else  
| 长度超限，退出  
end
```

4.3.1 字符串的顺序存储实现

```
void StrConcat(String s, String t)
{
    int n, m;
    Position i;
    n = s->length;
    m = t->length;
    if ((n+m)<kMaxSize) { /* 若串接后不超长 */
        for (i=0; i<=m; i++) { /* i=m处理结束符 */
            s->data[n+i] = t->data[i];
        }
        s->length = n + m;
    }
    else {
        printf("错误: 连接将导致字符串长度超限。\\n");
    }
}
```

4.3.1 字符串的顺序存储实现

■ 字符串比较

- 将这两个字符串从左到右逐个字符按照其ASCII码值进行比较。
- 如果两个字符串长度相等，且每一个相应位置上的字符都相同，则这两个字符串相等，如"abc"与"abc"相等。
- 如果两个字符串长度不相等，但较短字符串所有对应位置上的字符都与较长字符串相同，则字符串长度长的字符串大于长度短的字符串。如"abc"<"abcdef"。
- 如果两个字符串在某一相应位置上的字符不相同，则以第一个不相同的位置上的字符比较结果作为两个字符串的比较结果。如"abh">"abf"。
- 算法4-5对两个字符串做比较，若 $s > t$ 返回+1；若 $s = t$ 返回0；若 $s < t$ 返回-1。

■ 时间复杂度：仅与s和t的最小长度有关，为 $O(\min(n, m))$

4.3.1 字符串的顺序存储实现

```
int StrCompare(String s, String t)
{
    int len, ret;
    Position i;
    len = min(s->length, t->length);
    i = 0;
    while ((i<len) && (s->data[i]==t->data[i])) {
        i++; /* 顺次比较s和t等长的部分 */
    }
    if (i==len) { /* 等长部分都相等 */
        if (s->length > len) {
            ret = 1;
        }
        else if (t->length > len) {
            ret = -1;
        }
        else { /* s = t */
            ret = 0;
        }
    } /* 否则等长部分不相等 */
    else if (s->data[i] > t->data[i]) {
        ret = 1;
    }
    else {
        ret = -1;
    }
    return ret;
}
```

算法4-5 字符串比较操作StrCompare(s,t)

输入：字符串s，字符串t
输出：s>t输出+1；s=t输出0；s<t输出-1
len ← Min(s.length, t.length)
i ← 0
while i<len 且 s.data[i] = t.data[i] do
 i ← i+1
end
if i=len then
 if s.length > len then
 ret ← 1
 else if t.length > len then
 ret ← -1
 else //s=t
 ret ← 0
 end
else if s.data[i] > t.data[i] then
 ret ← 1
else
 ret ← -1
end
return ret

4.3.2 字符串的链接存储实现

- 将每个结点存放一个字符的字符串链接存储方式称为**非紧缩链接存储结构**，一个结点存放多个字符的字符串链接存储方式称为**紧缩链接存储结构**，也称为**块链存储结构**。

- **例：字符串"abcdef"**

- 非紧缩链接存储结构：



- 紧缩链接存储结构：



- 为容易理解，下面对非紧缩链接存储结构的字符串的基本操作进行介绍

4.3.2 字符串的链接存储实现

```
typedef struct StringNode *Position; /* 指针即结点位置 */
struct StringNode {
    char data; /* 存储数据 */
    Position next; /* 链接存储中下一个元素的位置 */
};
typedef struct StringHeadNode *String;
struct StringHeadNode {
    Position head; /* 字符串头指针, 初始化为NULL */
    int length; /* 字符串长度 */
};

String InitStr()
{ /* 初始化一个空的字符串 */
    String s;
    s = (String)malloc(sizeof(struct StringHeadNode));
    s->head = NULL;
    s->length = 0;
    return s;
}
```

4.3.2 字符串的链接存储实现

■ 字符串插入

相对于顺序存储结构，使用链接存储结构的字符串插入操作较为简单，无需担心字符串长度超限的问题

- 需要找到字符串s中位于pos位置的链表元素
- 需要找到插入字符串t的末尾

■ 时间复杂度

该操作时间复杂度是 $O(n+m)$ ，其中n和m分别是两个字符串的长度

算法4-6 字符串插入操作StrInsert(s, pos, t)

输入：字符串s，要插入的位置 $pos \geq 1$ ，需要插入的字符串t

输出：完成插入后的字符串s。若不存在位置pos，则s不变

```
1      flag ← NormalCode
2      if t.length > 0 then //若t不是空串
3          | tail ← t.head
4          | while tail.next ≠ NIL do //找到t的最后一个元素
5              | | tail ← tail.next
6          | end
7          | if s.length > 0 then //若s不是空串
8              | | p ← s.head
9              | | count ← 1
10             | | while p ≠ NIL 且 count < pos-1 do //找第pos个元素的前一个元素
11                 | | | count ← count + 1
12                 | | | p ← p.next
13             | | end
14             | | if count = (pos-1) then //将t插在p的后面
15                 | | | tail.next ← p.next
16                 | | | p.next ← t.head
17             | | else if pos = 1 then //t插在s的头
18                 | | | tail.next ← s.head
19                 | | | s.head ← t.head
20             | | else
21                 | | | flag ← ErrorCode //输入错误：不存在位置pos
22             | | end
23         | else //若s是空串
24             | | s ← t
25         | end
26     end
27     if flag ≠ ErrorCode then //正常完成插入
28         | s.length ← s.length + t.length
29     end
```

4.3.2 字符串的链接存储实现

```
void StrInsert(String s, int pos, String t){
    int flag, count;
    Position p, tail;
    flag = NormalCode;
    if (t->length>0) { /* 若t不是空串 */
        tail = t->head;
        while (tail->next != NULL) { /* 找到t的最后一个元素 */ tail = tail->next;}
        if (s->length>0) { /* 若s不是空串 */
            p = s->head;
            count = 1;
            while ((p!=NULL) && (count<(pos-1))) { /* 找第pos个元素的前一个元素 */
                count++;
                p = p->next;}
            if (count==(pos-1)) { /* 将t插在p的后面 */
                tail->next = p->next;
                p->next = t->head;}
            else if (pos == 1) { /* t插在s的头 */
                tail->next = s->head;
                s->head = t->head;}
            else {
                printf("错误: 指定插入位置不存在。 \n");
                flag = ErrorCode;}}
        else { /* 若s是空串 */ s = t;}}
    if (flag != ErrorCode) { /* 正常完成插入 */
        s->length = s->length + t->length;}}
```

4.3.2 字符串的链接存储实现

■ 字符串删除

返回字符串s从pos位置开始删除长度为len的子串后的字符串

- 将pos-1位置上的链表元素的next指针指向原字符串s中的第pos+len个元素
- 释放被删除的结点空间

■ 时间复杂度

这个操作的时间复杂度与s的长度有关，为 $O(n)$

算法4-7 字符串删除操作StrRemove(s, pos, len)

输入：字符串s，要删除的位置 $pos \geq 1$ ，删除的字符个数len

输出：删除子串后的字符串s。若删除位置不存在，则s不变。

```
1.      if s.length>0 then //若s不是空串
2.      | p ← s.head
3.      | count ← 1
4.      | while p ≠ NIL 且 count<pos-1 do //找第pos个元素的前一个元素
5.      | | count ← count + 1
6.      | | p ← p.next
7.      | end
8.      | if pos=1 或 (p ≠ NIL 且 count=pos-1) then //将p后面的len个结点删除
9.      | | if pos=1 then
10.     | | | deleted ← s.head
11.     | | | else
12.     | | | | deleted ← p.next
13.     | | | end
14.     | | count ← 0
15.     | | while deleted ≠ NIL 且 count < len do //不足len个则一直删到末尾
16.     | | | t ← deleted.next
17.     | | | delete deleted
18.     | | | count ← count + 1
19.     | | | s.length ← s.length - 1
20.     | | | deleted ← t
21.     | | end
22.     | | if pos = 1 then
23.     | | | s.head ← deleted
24.     | | | else
25.     | | | | p.next ← deleted
26.     | | | end
27.     | end
28.     end
```

4.3.2 字符串的链接存储实现

```
void StrRemove(String s, int pos, int len){
    Position p, deleted, t;
    int count;
    if (s->length > 0) {
        p = s->head;
        count = 1;
        while ((p!=NULL) && (count<(pos-1))) { /* 找第pos个元素的前一个元素 */
            count++;
            p = p->next; }
        if ((pos==1) || (p!=NULL && count==(pos-1))) {
            /* 将p后面的len个结点删除 */
            if (pos==1) {
                deleted = s->head; }
            else {
                deleted = p->next; }
            count = 0;
            while ((deleted!=NULL) && (count<len)) { /* 不足len个则一直删到末尾 */
                t = deleted->next;
                free(deleted);
                count++;
                s->length--;
                deleted = t; }
            if (pos==1) { s->head = deleted; }
            else { p->next = deleted; }
        }
    }
}
```

4.3.2 字符串的链接存储实现

■ 字符串截取

返回字符串s从pos位置开始的长度为len的子串

- 构造新串

- 将字符串s从pos之后len个字符一一赋值给新串。

■ 时间复杂度

与顺序存储不同，这个操作必须首先找到截取的起始位置，时间复杂度就不仅与子串长度有关了，最坏时间复杂度为 $O(n)$ 。

算法4-8 字符串截取子串操作SubString(s, pos, len)

输入：字符串s，开始截取的位置 $\text{pos} \geq 1$ ，截取的字符个数len

输出：截取的字串

```
1.      InitStr(sub_s) //初始化新串sub_s
2.      if s.length>0 then //若s不是空串
3.      | p ← s.head
4.      | count ← 1
5.      | while p ≠ NIL 且 count<pos do //找第pos个元素
6.      | | count ← count + 1
7.      | | p ← p.next
8.      | end
9.      | if p ≠ NIL 且 count=pos then //将s串从pos之后的len个字符复制到sub_s
10.     | | count ← 0
11.     | | sub_s.head ← new StringNode() //创建临时空头结点
12.     | | tail ← sub_s.head
13.     | | while p ≠ NIL 且 count<len do //若从pos开始不到len个字符，就截取到s的末尾
14.     | | | t ← new StringNode(p.data, NIL) //复制一个新结点
15.     | | | tail.next ← t //将新结点接到sub_s的末尾
16.     | | | tail ← tail.next
17.     | | | sub_s.length ← sub_s.length + 1 //sub_s 长度加1
18.     | | | p ← p.next
19.     | | | count ← count + 1
20.     | | end
21.     | end
22.     end
23.     tail ← sub_s.head
24.     sub_s.head ← tail.next
25.     delete tail //删除临时空头结点
26.     return sub_s
```

4.3.2 字符串的链接存储实现

```
String SubString(String s, int pos, int len)
{
    String sub_s;
    Position p, t, tail;
    int count;
    sub_s = InitStr(); /* 初始化新串sub_s */
    if (s->length>0) { /* 若s不是空串 */
        p = s->head;
        count = 1;
        while ((p!=NULL) && (count<pos)) { /* 找第pos个元素 */
            count++;
            p = p->next;
        }
        if ((p!=NULL) && (count==pos)) {
            /* 将s串从pos之后的len个字符复制到sub_s */
            count = 0;
            /* 创建临时空头结点 */
            sub_s->head = (Position)malloc(sizeof(struct StringNode));
            tail = sub_s->head;
```


4.3.2 字符串的链接存储实现

```
while ((p!=NULL) && (count<len)) {
    /* 若从pos开始不到len个字符, 就截取到s的末尾 */
    t = (Position)malloc(sizeof(struct StringNode));
    t->data = p->data; /* 复制一个新结点 */
    t->next = NULL;
    tail->next = t; /* 将新结点接到sub_s的末尾 */
    tail = tail->next;
    sub_s->length++; /* sub_s 长度加1 */
    p = p->next;
    count++;
}
}
tail = sub_s->head;
sub_s->head = tail->next;
free(tail);
return sub_s;
}
```

4.3.2 字符串的链接存储实现

■ 字符串连接

将字符串t连接在字符串s末尾。
与顺序存储不同，这里不需要将字符串t中的字符复制到字符串s末尾，只需要找到字符串s末尾并将字符串t接在后面即可。

■ 时间复杂度

这个操作的时间复杂度与字符串t的长度无关，仅与字符串s的长度成正比，为 $O(n)$ 。

算法4-9 字符串连接操作StrConcat(s,t)

输入：字符串s，字符串t

输出：返回后面联接t而成的字符串s

```
1.      if s.length > 0 then //若s非空串
2.      | p ← s.head
3.      | while p.next ≠ NIL do //找到s的最后一个结点
4.      | | p ← p.next
5.      | end
6.      | p.next ← t.head
7.      else //若s是空串
8.      | s.head ← t.head
9.      end
10.     s.length ← s.length + t.length
```

4.3.2 字符串的链接存储实现

```
void StrConcat(String s, String t)
{
    Position p;
    if (s->length > 0) { /* 若s非空串 */
        p = s->head;
        while (p->next != NULL) {
            p = p->next; /* 找到s的最后一个结点 */
        }
        p->next = t->head;
    }
    else { /* 若s是空串 */
        s->head = t->head;
    }
    s->length += t->length;
}
```

4.3.2 字符串的链接存储实现

■ 字符串比较

返回字符串s和字符串t的大小关系。

- 若 $s > t$ 返回+1
- 若 $s = t$ 返回0
- 若 $s < t$ 返回-1

■ 时间复杂度

这个操作的时间复杂度也是仅与s和t的最小长度有关，为 $O(\min(n, m))$

算法4-10 字符串比较操作StrCompare(s,t)

输入：字符串s，字符串t

输出： $s > t$ 输出+1； $s = t$ 输出0； $s < t$ 输出-1

```
1.      sp ← s.head
2.      tp ← t.head
3.      while sp ≠ NIL 且 tp ≠ NIL 且 sp.data = tp.data do
4.      | sp ← sp.next
5.      | p ← tp.next
6.      end
7.      if sp ≠ NIL 且 tp = NIL then
8.      | ret ← 1
9.      else if sp = NIL 且 tp ≠ NIL then
10.     | ret ← -1
11.     else if sp = NIL 且 tp = NIL then //s=t
12.     | ret ← 0
13.     else if sp.data > tp.data then
14.     | | ret ← 1
15.     else //sp.data < tp.data
16.     | | ret ← -1
17.     end
18.     return ret
```

4.3.2 字符串的链接存储实现

```
int StrCompare(String s, String t)
{
    Position sp, tp;
    int ret;
    sp = s->head;
    tp = t->head;
    while ((sp!=NULL) && (tp!=NULL) && (sp->data==tp->data)) {
        /* 顺次比较s和t等长的部分 */
        sp = sp->next;
        tp = tp->next;
    }
    if ((sp!=NULL) && (tp==NULL)) {ret = 1;}
    else if ((sp==NULL) && (tp!=NULL)) {ret = -1;}
    else if ((sp==NULL) && (tp==NULL)) {ret = 0;}
    else if (sp->data > tp->data) {ret = 1;}
    else {ret = -1;}
    return ret;
}
```

- 4.1 问题引入：模式匹配
- 4.2 字符串的定义与结构
- 4.3 字符串的存储实现
- 4.4 字符串的模式匹配
- 4.5 拓展延伸
- 4.6 应用场景：基因测序

4.4 字符串的模式匹配

- **概念：**在字符串s中找出与字符串t相等的子串的操作称为字符串的模式匹配，又称为子串的定位操作。其中字符串s称为主串或目标串，字符串t称为模式串。
- **模式匹配算法：**朴素字符串匹配算法（BF算法）、KMP算法、BM算法、KR算法、Sunday算法等。
- **形式化描述：**假设目标串s使用一个长度为n的字符数组 $s[0,1,\dots,n-1]$ 表示，模式串t使用一个长度为m ($m \leq n$)的数组 $t[0,1,\dots,m-1]$ 表示，如果存在p ($0 \leq p \leq n - m$)，使得 $s[p+0, p+1, \dots, p+m-1] = t[0, 1, \dots, m-1]$ ，则p被称为一个有效位移。字符串匹配就是从字符串s中找出存在的有效位移p。

4.4.1 朴素模式匹配算法

- 朴素模式匹配算法是字符串模式匹配算法中最简单的暴力解法，又称为BF（Brute Force）算法。
 - 实现：枚举每个目标串s与模式串t等长的子串，判断是否匹配
 - 将模式串t的第0位字符与目标串s的第0位字符对齐，然后依次比对每个字符，
 - 都相等，则匹配成功
 - 若s和t对应位置上的字符不相等，则将s整体后移1位重新从模式串t的第0位开始依次比对。

算法4-11 朴素字符串匹配算法PatternMatchBF(s, t)

输入：目标串s与模式串t

输出：返回首个有效位移p，匹配失败则返回NIL

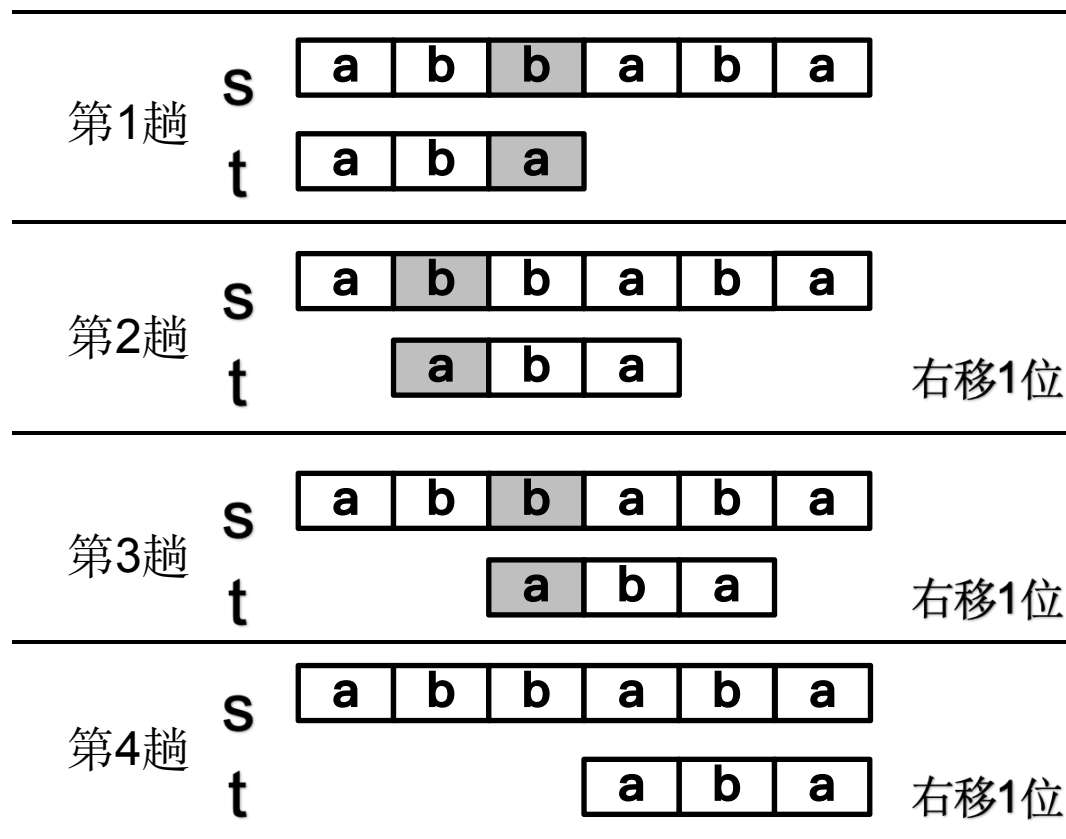
```
n ← s.length
m ← t.length
for p ← 0 to n-m do
| for i ← 0 to m-1 do
| | if s.data[p+i] ≠ t.data[i] then
| | | break
| | end
| end
| if i=m then
| | break
| end
end
if p > n-m then
| p ← NIL
end
return p
```


4.4.1 朴素模式匹配算法

```
Position PatternMatchBF(String s, String t)
{
    int n, m, i;
    Position p;
    n = s->length;
    m = t->length;
    for (p=0; p<=(n-m); p++) {
        for (i=0; i<m; i++) {
            if (s->data[p+i] != t->data[i]) {
                break; /* 不匹配则从下一个位置p开始 */
            }
        }
        if (i==m) { /* 匹配成功 */
            break; }
    }
    if (p>(n-m)) { /* 匹配不成功 */
        p = NIL; }
    return p;
}
```

4.4.1 朴素模式匹配算法

- 例：目标串s="abbaba", 模式串t="aba"



4.4.1 朴素模式匹配算法

■ 朴素模式匹配算法时间复杂度分析

- 最好情况下，仅需匹配 m 次，此种情况时间复杂度为 $O(m)$
- 最坏情况下，需要移动 $n-m+1$ 次，每次匹配 m 次，时间复杂度为 $O(nm-m^2)$ ，即 $O(nm)$
- 在实际运行过程中字符串匹配情况复杂多变，朴素模式匹配算法的执行时间通常取上界 $O(nm)$

4.4.1 朴素模式匹配算法

- 字符串匹配算法是其它部分字符串操作的基础，例如字符串替换操作需要先使用匹配操作找到位置，接着进行字符串的删除和插入替换

算法 4-12 字符串替换算法Replace(s, sub_s, t)

输入：字符串s，被替换的子串sub_s，替换目标字符串t

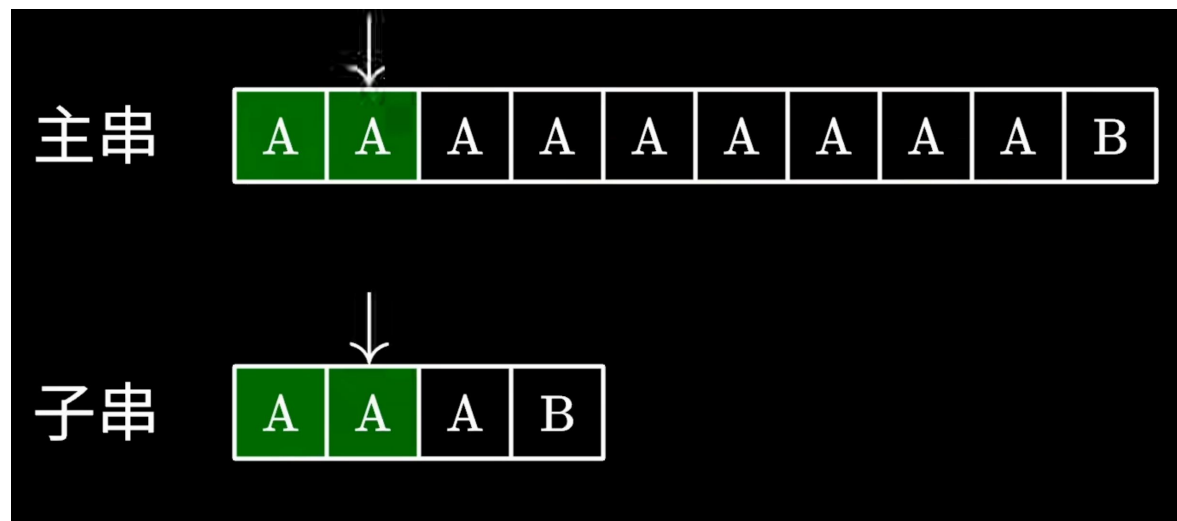
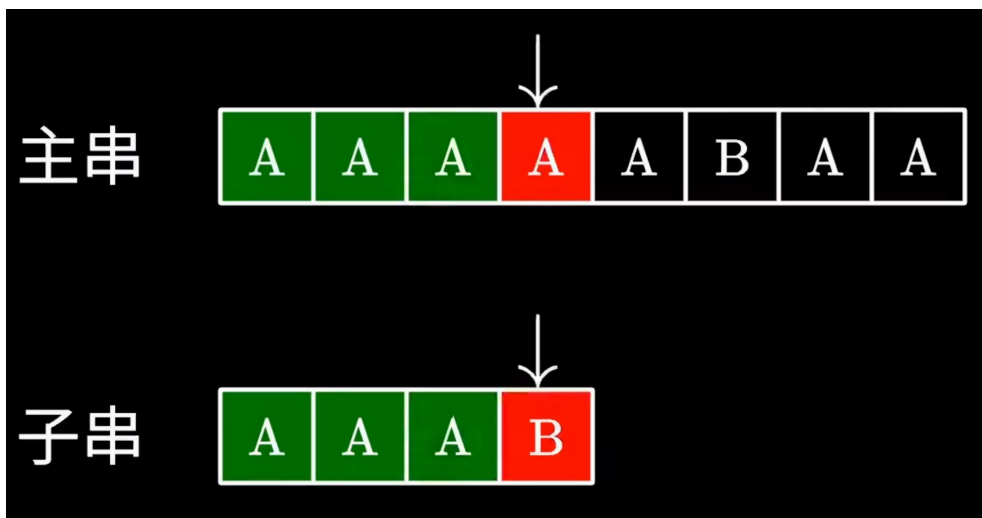
输出：替换字符串s中的所有子串sub_s为字符串t后的字符串s

```
1  len ← sub_s.length
2  m ← t.length
3  pos ← 0
4  while pos ≠ NIL do
5    | pos ← PatternMatchBF(s, sub_s) //从s中找到第一次出现的sub_s
6    | if pos ≠ NIL then
7      | | StrRemove(s, pos+1, len) //删除sub_s
8      | | StrInsert(s, pos+1, t) //插入t
9      | | s.length ← s.length - len + m
10   | end
11 end
```

4.4.1 朴素模式匹配算法

```
void Replace(String s, String sub_s, String t)
{
    Position pos;
    int len, m;
    len = sub_s->length;
    m = t->length;
    pos = 0;
    while (pos != NIL) {
        pos = PatternMatchBF(s, sub_s); /* 从s中找到第一次出现的sub_s */
        if (pos != NIL) {
            StrRemove(s, pos+1, len); /* 删除sub_s */
            StrInsert(s, pos+1, t); /* 插入t */
        }
    }
}
```

4.4.1 朴素模式匹配算法



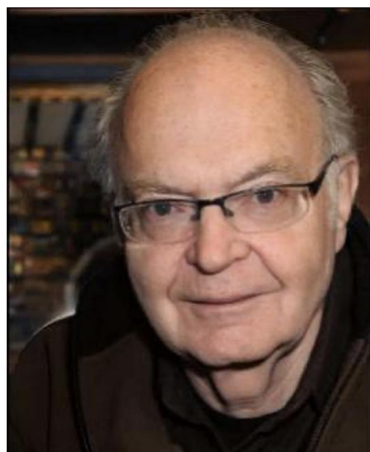
■ 匹配失败时

- 主串（目标串）：回溯到目标串匹配的下一个位置， $p++$
- 子串（模式串）：回溯到模式串起点， $i=0$

■ 能否避免主串回溯，从 $O(nm)$ 减小为线性 $O(n)$ 复杂度？

4.4.2 KMP 算法

- **提出：** KMP算法由D.E.**K**nuth, J.H.**M**orris和V.R.**P**ratt提出的



克鲁特**K**nuth
《计算机程序设计艺术》



莫里斯**M**orris

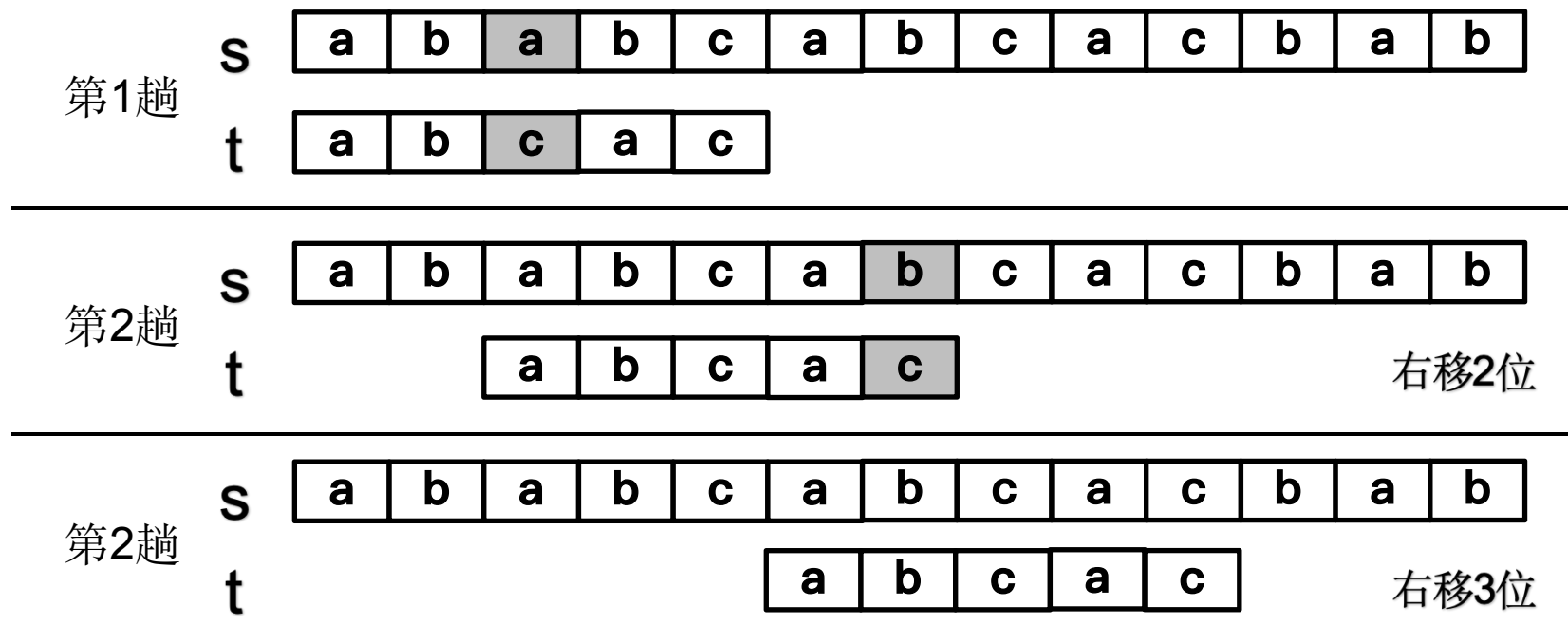


普拉特**P**ratt

- **主要思想：** 假设BF算法在失配时已经匹配到了字符串的第 j 位，则说明目标串与模式串的前 $j-1$ 位是匹配成功的，利用已匹配的信息可对BF算法进行优化。

4.4.2 KMP 算法

- 例：目标串 $s = \text{"ababcabcacbab"}$ ，模式串 $t = \text{"abcac"}$



- 可以看出在比较的时候，模式串向后移动了更多的位数，因此能够更快地完成模式匹配

4.4.2 KMP 算法

- 思考：当出现失配情况时，如何将模式串向右移动正确的位数？
- 字符串相关概念：
 - 前缀：从长度为 n 的字符串第0位开始，第 i 位 ($0 \leq i < n-1$) 结束的任意子串。对字符串 s ，其前缀可表示为 $s[0...i]$ ($0 \leq i < n-1$)。字符串的所有前缀构成的集合称为前缀集合。
 - 后缀：从长度为 n 的字符串第 i 位 ($0 < i \leq n-1$) 开始，最后一位结束的任意子串。对字符串 s ，其后缀可表示为 $s[i...n-1]$ ($0 < i \leq n-1$)。字符串的所有后缀构成的集合称为后缀集合。
 - 公共前后缀：字符串的前缀集合与后缀集合中相同的子串。
 - 最长公共前后缀：字符串的前缀集合与后缀集合中相同的长度最长的子串。

4.4.2 KMP 算法

■ 例：字符串"aabaa"

- 前缀集合为{"a","aa","aab","aaba"}
- 后缀集合为{"a","aa","baa","abaa"}
- 其公共前后缀包括"a"和"aa"两个子串
- 最长公共前缀为"aa"这个子串



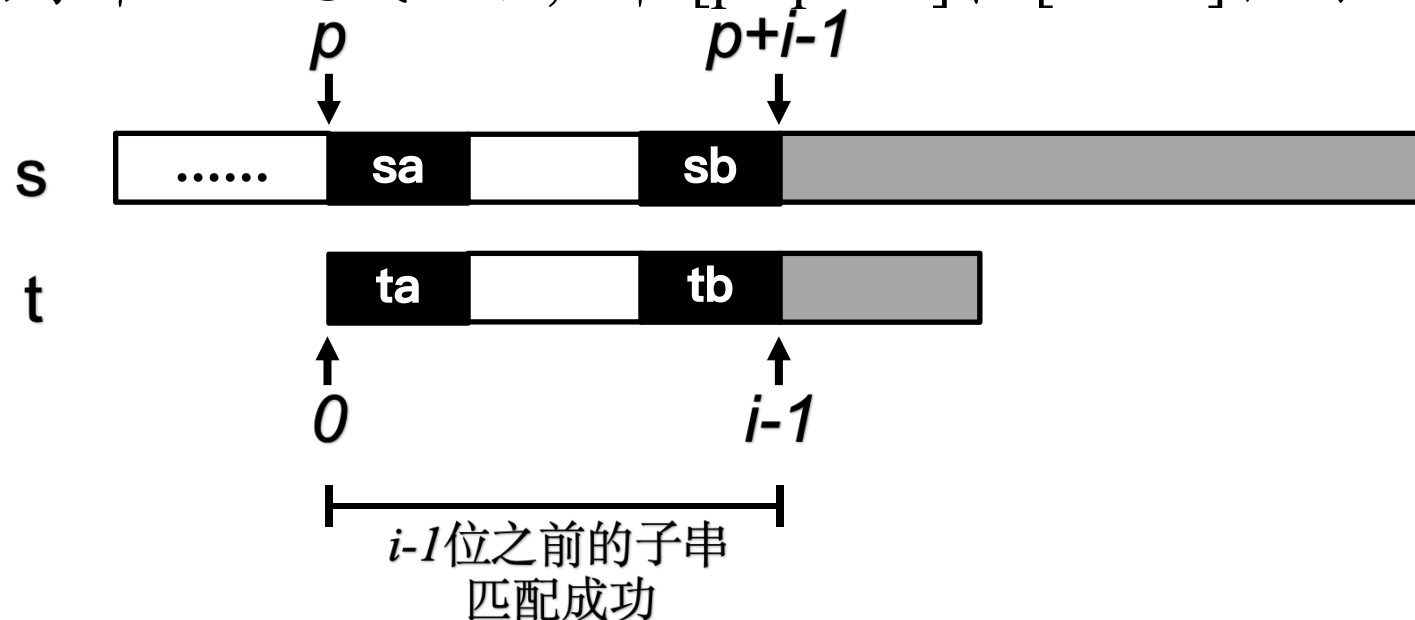
- 目标串的最长公共前后缀 $sa=sb$



- 模式串的最长公共前后缀 $ta=tb$

4.4.2 KMP 算法

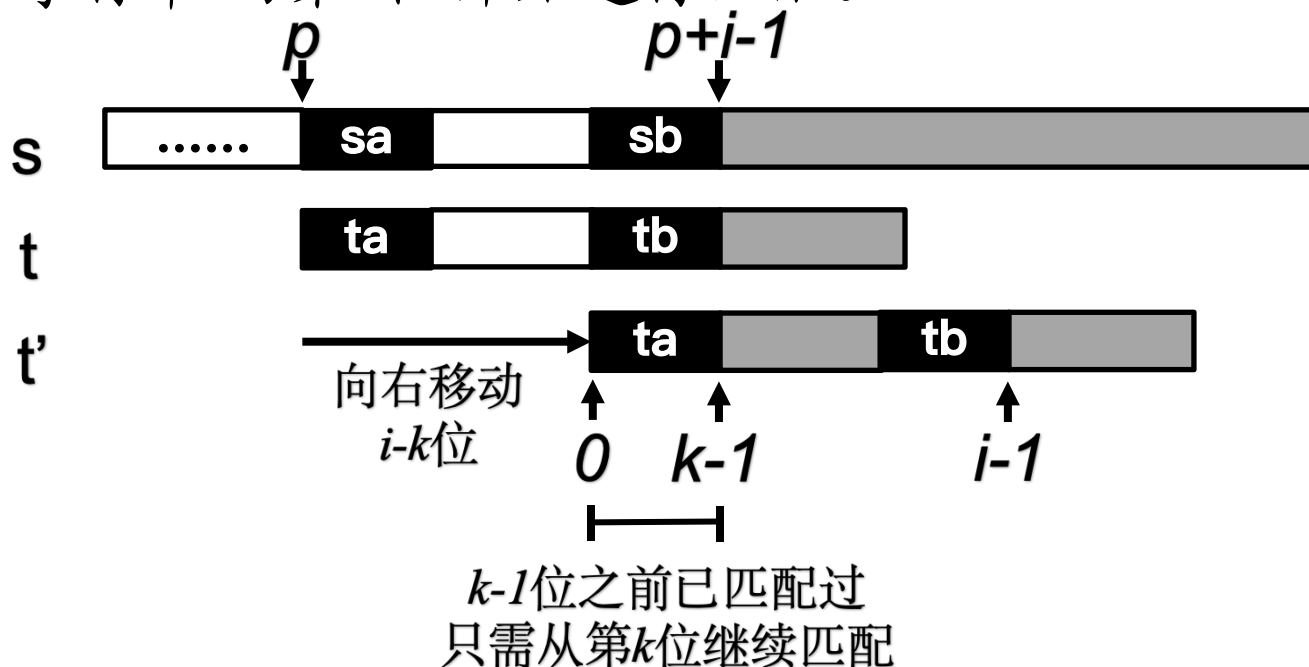
- 假设当前模式串 t 已经移动到目标串 s 的某一位置 p ，即字符串 t 的第 0 位字符与 s 的第 p 位字符对齐，正在对模式串 t 的第 i 位进行匹配，则说明模式串 t 的 $i-1$ 位之前的子串已经完成匹配，即 $s[p \dots p+i-1]$ 和 $t[0 \dots i-1]$ 相同。



- 由于 $s[p \dots p+i-1]$ 和 $t[0 \dots i-1]$ 已成功匹配，所以显然有 $sa=tb=ta=tb$ 成立，记最长公共前后缀长度为 k 。

4.4.2 KMP 算法

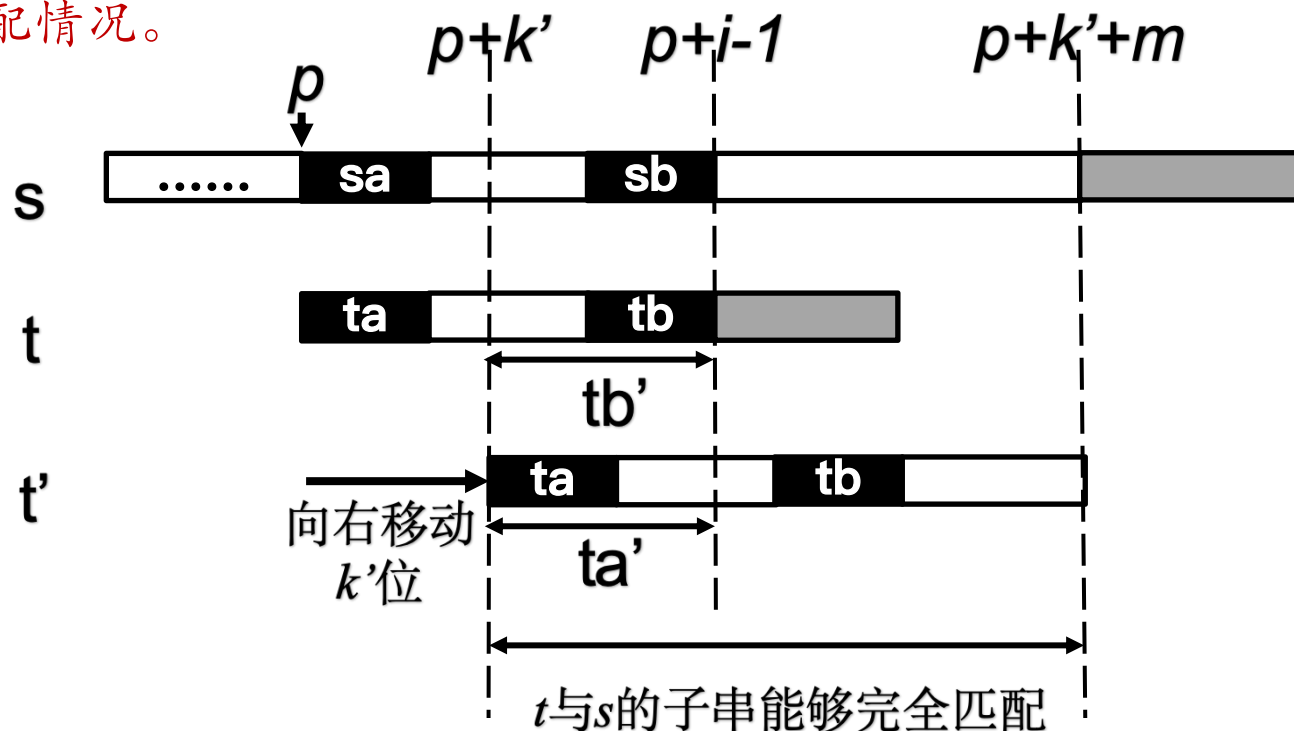
- 当模式串 t 的第 i 位与目标串 s 的第 $p+i$ 位出现失配情况时，如果把字符串 t 向右移动 $i-k$ 位，使 sb 和 ta 对齐，如图4-7所示。那么，由于 $sb=ta$ ，所以只需要从 ta 的下一位即字符串 t 的第 k 位与字符串 s 的第 $p+i$ 位再开始进行匹配即可，就不需要再从字符串 t 的第0位开始进行匹配。



这种方法虽然能够大大提升匹配的效率，但多位移动中会不会漏掉一些匹配的情况？

4.4.2 KMP 算法

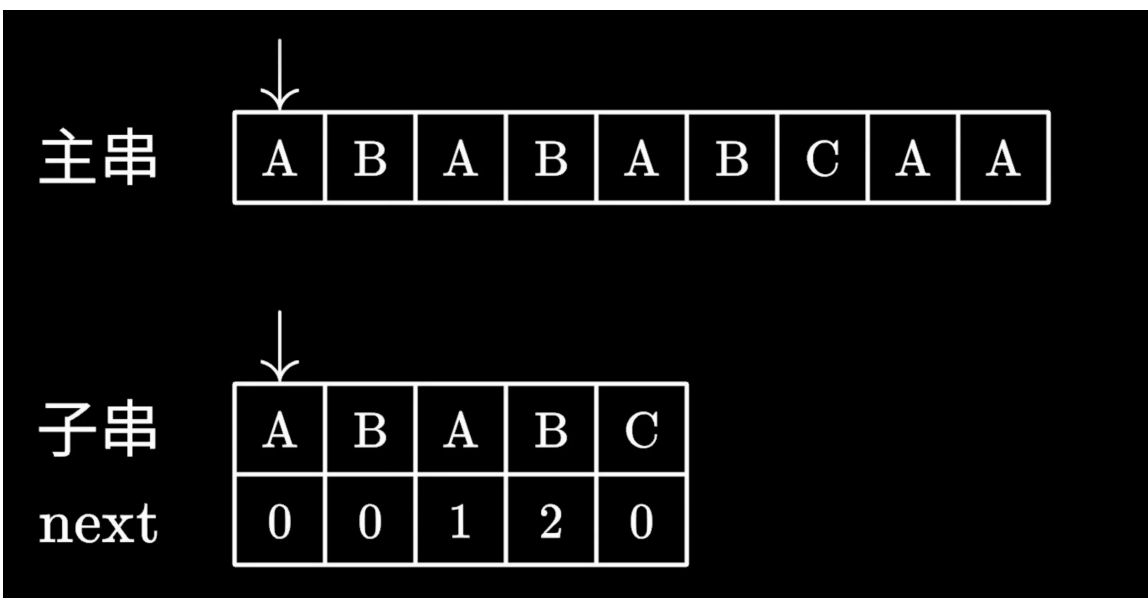
- 记模式串 t 的长度为 m ，假设出现了漏掉情况
 - s 的子串 $s[p+k' \dots p+k'+m]$ 与字符串 t 完全相同，那么子串 $s[p+k' \dots p+i-1]$ 与模式串 t 的前缀 ta' 相同
 - 由于前 $i-1$ 位已经成功匹配，故子串 $s[p+k' \dots p+i-1]$ 与模式串 t 的后缀 tb' 也相同
 - 模式串 t 的前缀 ta' 与后缀 tb' 也相同，是 t 的公共前后缀
 - 由于 ta' 长度大于 ta ，与 ta 是 t 的最长公共前后缀矛盾，所以KMP算法在多位移动中不会漏掉可能出现的匹配情况。



4.4.2 KMP 算法

■ 模式串移动多少位?

- 前缀与后缀对齐
- next数组放最大公共前后缀长度（或者最大公共前后缀长度-1）
- 注：教材next数组 = 最大公共前后缀长度-1



4.4.2 KMP 算法

KMP算法中向右移动的位数由已匹配部分字符串的最长公共前后缀的长度决定。
如何计算这个位数？

字符串特征向量：长度为m的字符串t的特征向量是一个m维向量，通常记作next，且用数组形式进行存储，所以特征向量也可非正式地称为**next数组**。next[i]表示字符串特征向量的第i位分量（ $0 \leq i < m$ ），其形式化定义为：

$$\text{next}[i] = \begin{cases} \text{满足 } t[0 \dots k] = t[i - k \dots i] \text{ 的最大 } k (\text{存在 } k, \text{ 且 } k < i) \\ -1, & \text{如果这样的 } k \text{ 不存在} \end{cases}$$

next[i]表示的真实含义是字符串t的子串t[0...i]的最长公共前后缀中的前缀最末尾的字符的位置。由于字符串从0位开始，所以t[0...i]的最长公共前后缀长度为next[i]+1。如果t[0...i]的最长公共前后缀不存在，那么将next[i]置为-1。

4.4.2 KMP 算法

- 例：字符串"abccac"的特征向量为 $[-1, -1, -1, 0, -1]$

i	子串 $t[0...i]$	前缀集合	后缀集合	最长 公共前后缀	特征向量 第 <i>i</i> 位分量
$i=0$	a	空	空	空	-1
$i=1$	ab	a	b	空	-1
$i=2$	abc	a,ab	c,bc	空	-1
$i=3$	abca	a,ab,abc	a,ca,bca	a	0
$i=4$	abccac	a,ab,abc,abca	c,ac,cac,bcac	空	-1

4.4.2 KMP 算法

■ KMP 算法流程

在next数组的基础上，在目标串与模式串在模式串的第i位失配时

- 将模式串右移，从模式串的第 $\text{next}[i-1]+1$ 位开始进行匹配，并不断执行以上步骤直到字符串末尾。

■ 时间复杂度

算法执行过程中目标串只向右移动，比较的复杂度为 $O(n)$ 。

算法4-14 字符串匹配的KMP算法PatternMatchKMP(s, t)

输入：目标串s，模式串t

输出：返回首个有效匹配位置p，匹配失败则返回NIL

```
n ← s.length
m ← t.length
p ← NIL
if n ≥ m then
  | GetNext (t, next)
  | i ← 0
  | j ← 0
  | while j < n 且 i < m do
  | | if s.data[j] = t.data[i] then
  | | | i ← i + 1
  | | | j ← j + 1
  | | else if i > 0 then
  | | | i ← next[i - 1] + 1
  | | else
  | | | j ← j + 1
  | | end
  | end
  | if i = m then
  | | p ← j - m
  | end
end
return p
```

4.4.2 KMP 算法

```
Position PatternMatchKMP(String s, String t)
{
    int n, m, i, j;
    Position p;
    int *next;
    n = s->length;
    m = t->length;
    p = NIL;
    if (n >= m) {
        next = (int *)malloc(sizeof(int)*m);
        GetNext(t, next); /* 计算next数组 */
        i = j = 0;
```

4.4.2 KMP 算法

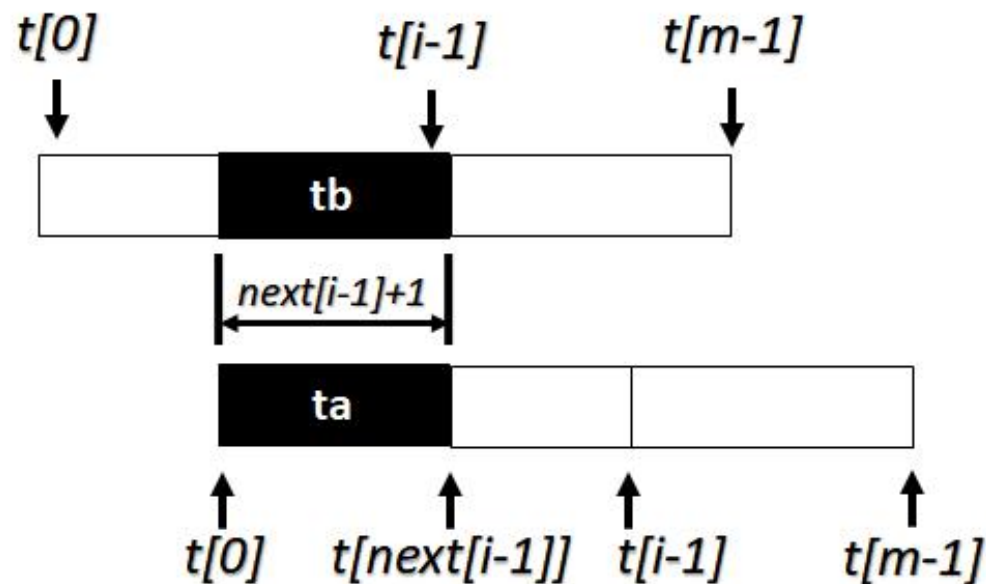
```
while ((j<n) && (i<m)) {  
    /* 若匹配成功则两指针同时向右移动 */  
    if (s->data[j] == t->data[i]) {  
        i++;  
        j++;  
    }  
    else if (i>0) { /* 失配且t指针未回退到左端 */  
        i = next[i-1] + 1; /* 计算t的匹配起始位置 */  
    }  
    else { /* 失配且t从头开始与s的下一段比较 */  
        j++;  
    }  
}  
if (i==m) { /* t完成了匹配 */  
    p = j - m; /* 匹配成功的位置 */  
}  
}  
return p;  
}
```

主串中的
指针 i 永不回退
时间复杂度 $O(n)$

4.4.2 KMP 算法

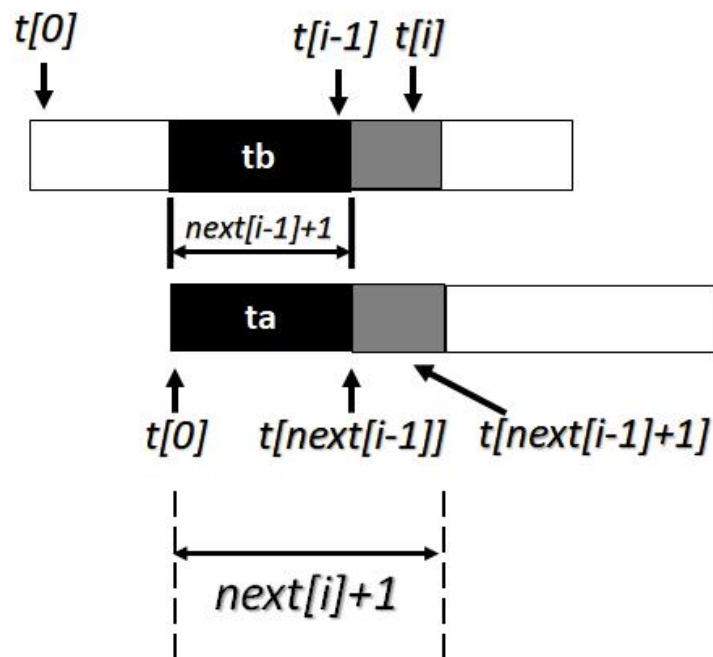
■ 求解next数组

- 若朴素计算next数组，时间复杂度为 $O(m^3)$ ，不能满足进行快速模式匹配的要求
- 一种时间复杂度为 $O(m)$ 的字符串特征向量算法：图示为在字符串 t 中 $\text{next}[i-1]$ 的含义，黑色背景部分为子串 $t[0\dots i-1]$ 的最长公共前后缀，前缀用 ta 表示，后缀用 tb 表示，有 $ta=tb$ 成立， $\text{next}[i-1]$ 表示的是 ta 最末位字符的下标， $\text{next}[i-1]+1$ 表示的是最长公共前后缀的长度。



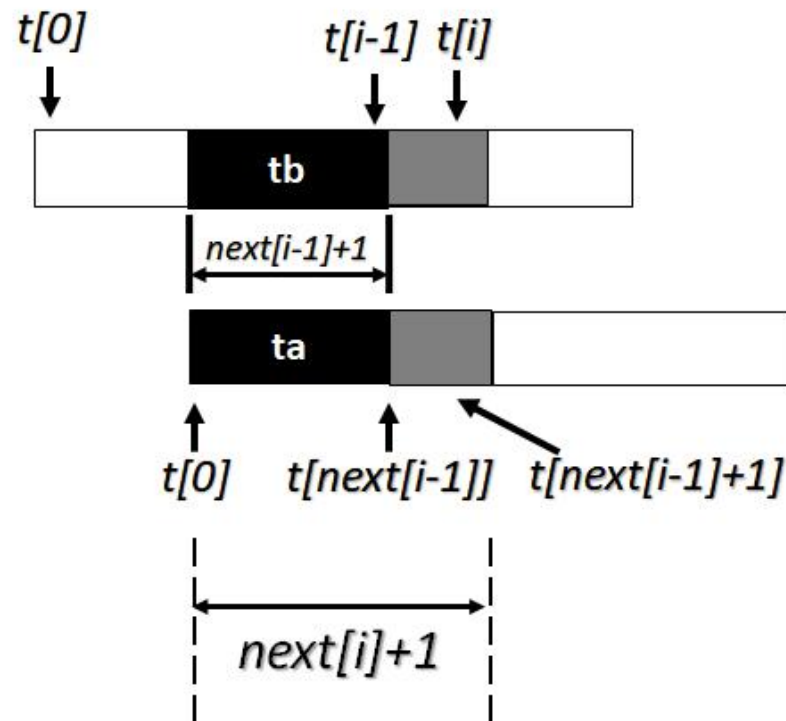
4.4.2 KMP 算法

- 考察tb后面一个字符，即t串第i位置上的字符 $t[i]$ 与ta后面一个字符，即字符串t第 $\text{next}[i-1]+1$ 位字符 $t[\text{next}[i-1]+1]$ 进行比较，则可对 $\text{next}[i]$ 的值进行推导。
- 如果 $t[i]$ 与 $t[\text{next}[i-1]+1]$ 相等，则tb加上字符 $t[i]$ 构成的子串与ta加上字符 $t[\text{next}[i-1]+1]$ 构成的子串能够完全匹配，这两个子串即为 $t[0\dots i]$ 的最长公共前后缀，故可推出 $\text{next}[i]=\text{next}[i-1]+1$



4.4.2 KMP 算法

- 如果当前 $t[i]$ 与 $t[\text{next}[i-1]+1]$ 不同，则说明 t 与其后一个字符形成的子串 $t[0\dots i]$ 的后缀，与 t 与其后一个字符形成的子串 $t[0\dots i]$ 的前缀并不匹配。
- 所以，在求解子串 $t[0\dots i]$ 的最长公共前缀时需要进一步缩小搜索范围，在子串 $t[0\dots \text{next}[i-1]]$ 中寻找可能的最长公共前后缀，即为 $t[0\dots i]$ 的最长公共前后缀，而子串 $t[0\dots \text{next}[i-1]]$ 的最长公共前后缀的前缀的末位下标即为 $\text{next}[\text{next}[i-1]+1]$ 。在此种情况下，有 $\text{next}[i] = \text{next}[\text{next}[i-1]+1]$ 成立。



4.4.2 KMP 算法

- 设 $j = \text{next}[i-1]$ ，则问题变为求子串 $t[0..j]$ 的最长公共前后缀长度 $\text{next}[j+1]$
 - 如果 $t[j]$ 与 $t[\text{next}[j-1]+1]$ 相等，则 $\text{next}[j] = \text{next}[j-1] + 1$
 - 否则， $\text{next}[j] = \text{next}[\text{next}[j-1] + 1]$
 - 设 j' 为 $\text{next}[j-1]$ ，继续按同样的方法求 $\text{next}[j']$ ，直到 $t[j']$ 与 $t[\text{next}[j'-1]+1]$ 相等或 j' 不合法即小于 0 为止。
- 计算 next 数组的复杂度为 $O(m)$ ，其中 m 为模式串的长度，故 KMP 算法的时间复杂度为 $O(n+m)$

算法4-13 求解字符串 t 的 next 数组 $\text{GetNext}(t, \text{next})$

```
输入：字符串  $t$   
输出：字符串  $t$  的  $\text{next}$  数组  
 $m \leftarrow t.\text{length}$   
 $\text{next}[0] \leftarrow -1$   
for  $i \leftarrow 1$  to  $m-1$  do // 求出  $\text{next}[1] \sim \text{next}[m-1]$   
|  $j \leftarrow \text{next}[i-1]$   
| while  $j \geq 0$  且  $t.\text{data}[i] \neq t.\text{data}[j+1]$  do  
| |  $j \leftarrow \text{next}[j]$   
| end  
| if  $t.\text{data}[i] = t.\text{data}[j+1]$  then  
| |  $\text{next}[i] \leftarrow j+1$   
| else  
| |  $\text{next}[i] \leftarrow -1$   
| end  
end
```

4.4.2 KMP 算法

```
void GetNext(String t, int next[])
{
    int m, i, j;
    m = t->length;
    next[0] = -1;
    for (i=1; i<m; i++) { /* 求出next[1]~next[m-1] */
        j = next[i-1];
        while ((j>=0) && (t->data[i]!=t->data[j+1])) {
            j = next[j];
        }
        if (t->data[i] == t->data[j+1]) {
            next[i] = j + 1;
        }
        else {
            next[i] = -1;
        }
    }
}
```


- 4.1 问题引入：模式匹配
- 4.2 字符串的定义与结构
- 4.3 字符串的存储实现
- 4.4 字符串的模式匹配
- 4.5 拓展延伸
- 4.6 应用场景：基因测序

4.5.1 正则表达式

- **概念：**正则表达式是对字符串的一种形式化描述方式，它通过组合一些具有特殊含义的字符，来实现对特定模式的字符串的描述。可以用来检索满足一定格式的字符串。
- **结构：**正则表达式由普通字符和元字符组成，普通字符包括了英文字母和数字等非元字符的字符，而元字符则具有特殊的含义。
 - 例：元字符“*”表示匹配前面的表达式任意次，“ab*”可匹配“a”，“ab”，“abb”等
- **实现：**对正则表达式建立有穷自动机，让此自动机仅接收满足规则的字符串，再让需要匹配的字符串在自动机上运行即可。

4.5.2 带有通配符的字符串匹配

- **概念：**通配符是一种特殊语句，主要有星号(*)和问号(?)，用来模糊搜索。带有通配符的字符串是一种带有星号和问号的字符串，主要用于字符串的模糊匹配。其中， '*'可以匹配0个或多个连续的任意字符， '?'可以匹配1个任意字符。
 - 例：模式串t为"a*b?c"，字符串"axyzbdc"、"abvc"均能够匹配模式串t。
- **注意：**带有通配符的字符串匹配定义与4.4节中提到的字符串模式匹配不同，其定义为判断目标串s和模式串t是否可以完全匹配，而非子串匹配。
 - 例如，模式串t可以和目标串s的子串匹配，但不能与目标串s完全匹配，则认为模式串t与目标串s失配。

4.5.2 带有通配符的字符串匹配

■ 一种递归算法:

- 1. s或者t其中一个已经到末尾了, 那么如果t的剩余字符都是 '*', 返回匹配, 否则返回不匹配;
- 2. 如果s的当前字符和t的当前字符相等, 继续向后移动;
- 3. 当s的当前字符和t的当前字符不相等, 分为三种情况:
 - (1) t的当前字符不是 '*' 或 '?', 返回不匹配;
 - (2) 若t的当前字符是 '?', 继续向后移动;
 - (3) t的当前字符是 '*', 那么可跳过s的0到多个字符, 再递归判断是否匹配

■ 时间复杂度

- 由于对每次递归判断过程, 目标串s和模式串t的指针i和j都是单调递增的, 所以复杂度为 $O(n+m)$ 。由于最多递归 $O(n)$ 次, 所以该算法的时间复杂度为 $O(n(n+m))$ 。

- 4.1 问题引入：模式匹配
- 4.2 字符串的定义与结构
- 4.3 字符串的存储实现
- 4.4 字符串的模式匹配
- 4.5 拓展延伸
- 4.6 应用场景：基因测序

4.6 应用场景：基因测序

- **示例场景：**生物信息学中已知一串复杂的病毒RNA序列，其规模可以达到百万级别。致病性的RNA序列段规模达到万级，此时科学家需要判断这串病毒的RNA序列中是否包含致病序列段。



- **解决方法：**将复杂的病毒RNA序列认为是目标串 s ，致病性的RNA序列段认为是模式串 t ，使用KMP算法解决模式匹配问题，判断模式串 t 是否为目标串 s 的子串，如果是，则可以认定该病毒具有致病性，反之则认为该病毒没有致病性。

小结

- 字符串和字符类型的线性表的异同点：
 - 字符串：主要是对字符串的整体操作
 - 线性表：是对表中元素的操作。
- 字符串的主要操作：
 - 字符串插入、字符串删除、字符串截取、字符串连接、字符串比较
- 字符串的两种实现方法：
 - 顺序存储结构实现、链接存储结构实现
- 介绍了字符串匹配中的几个重要算法
 - 朴素算法、**KMP算法**、BM算法等。

小结

- 串是内容受限的线性表，它限定表中的元素为字符。
- 常见的模式匹配，包括BF匹配算法和KMP匹配算法
- BF暴力模式匹配算法时间复杂度 $O(m*n)$ ，KMP匹配算法时间复杂度 $O(m+n)$
- BF暴力模式匹配算法失配时，主串 i 回溯到 $i-j+2$ ，模式串 j 回溯到 $j=1$
- KMP匹配算法失配时，主串 i 不需要回溯，模式串 j 回溯与模式串有关，使用 $next[j]$ 表示



字符串常见操作

操作名称	C 语言 (char* /char [])	C++ (std::string)	Java (String)	Python (str)
StrInsert(s, pos, t)	无直接函数, 需手动实现 (内存分配 + 拼接)	s.insert(pos, t)	s.substring(0, pos) + t + s.substring(pos)	s[:pos] + t + s[pos:] (切片拼接)
StrRemove(s, pos, len)	无直接函数, 需手动覆盖 + 补终止符	s.erase(pos, len)	s.substring(0, pos) + s.substring(pos+len)	s[:pos] + s[pos+len:] (切片拼接)
SubString(s, pos, len)	strncpy (需手动分配内存 + 补 '\0')	s.substr(pos, len)	s.substring(pos, pos+len)	s[pos:pos+len] (切片)
StrLength(s)	strlen(s)	s.size() / s.length()	s.length()	len(s)
StrConcat(s, t)	strcat /strncat (需确保 s 内存足够)	s + t / s.append(t)	s + t / s.concat(t)	s + t / ''.join([s,t])
StrCompare(s, t)	strcmp (返回正 / 0 / 负数)	s.compare (t) (返回正 / 0 / 负数)	s.compareTo (t) (返回正 / 0 / 负数)	自定义: (s>t)-(s<t) (返回 1/0/-1)
PatternMatch(s, t)	strstr (返回指针, 未找到为 NULL)	s.find (t) (未找到返回npos)	s.indexOf (t) (未找到返回 - 1)	s.find (t) (未找到返回 - 1)
Replace(s, sub_s, t)	strstr + 替换	循环配合 s.find () + s.replace ()	s.replace (sub_s, t) (全局替换)	s.replace (sub_s, t) (全局替换)



字符串高频必刷题 (LeetCode)

分类	题号	题目	考察点	难度
基础操作	344	Reverse String	双指针反转字符串	Easy
	415	Add Strings	模拟字符串加法	Easy
	14	Longest Common Prefix	字符串前缀比较	Easy
	58	Length of Last Word	字符串遍历与切分	Easy
双指针	125	Valid Palindrome	双指针判断回文字符串	Easy
	5	Longest Palindromic Substring	双指针中心扩展求最长回文子串	Medium
	151	Reverse Words in a String	翻转单词顺序，空格处理	Medium
哈希计数	242	Valid Anagram	哈希表统计字符频次	Easy
	49	Group Anagrams	哈希表分组异位词	Medium
滑动窗口	3	Longest Substring Without Repeating Characters	滑动窗口去重，最长无重复子串	Medium
	76	Minimum Window Substring	最小覆盖子串，滑动窗口经典题	Hard
	567	Permutation in String	字符计数匹配，滑动窗口子串判断	Medium



字符串高频必刷题 (LeetCode)

分类	题号	题目	考察点	难度
栈混合应用	394	Decode String	栈模拟字符串解码（嵌套括号结构）	Medium
	20	Valid Parentheses	括号匹配，字符串栈应用（与栈章节交叉）	Easy
字符匹配	28	Find the Index of the First Occurrence in a String	子串匹配，KMP/BF算法	Easy
	6	Zigzag Conversion	按行构造字符串	Medium
数字与编码	13	Roman to Integer	模拟罗马数字解析	Easy
	12	Integer to Roman	数值转罗马数字	Medium
	8	String to Integer (atoi)	字符串解析整数，状态机实现	Medium
表达式解析	224	Basic Calculator	表达式解析与括号处理	Hard
动态规划	72	Edit Distance	编辑距离，字符串DP经典题	Hard
	10	Regular Expression Matching	正则匹配DP，复杂状态转移	Hard
	97	Interleaving String	交错字符串判断	Medium