

2.4 String



- Definition and operations
- Applications

2.4.1. Definition



- **串**: 零个或多个**字符**组成的有限**序列**。
- **串长度**: 串中所包含的字符个数。
- **空串**: 长度为0的串, 记为: “”。
- **非空串**通常记为: $S = "s_1 s_2 \dots s_n"$
 - 其中: S 是串名, 双引号是**定界符**, 双引号引起来的部分是串值, $s_i (1 \leq i \leq n)$ 是一个任意字符。
 - 字符集: **ASCII**码、扩展**ASCII**码、**Unicode**字符集
- **子串**: 串中任意个连续的字符组成的子序列。
- **主串**: 包含子串串。
- **子串的位置**: 子串的第一个字符在主串中的序号。

ADT String{



数据对象: $D = \{ a_i | a_i \in \text{CharacterSet}, i=1,2,\dots,n, n \geq 0 \}$

数据关系: $R = \{ \langle a_{i-1}, a_i \rangle | a_{i-1}, a_i \in D, i=2,3,\dots,n \}$

基本操作:

StrAssign(t, chars)

初始条件: chars是一个字符串常量。

操作结果: 生成一个值为chars的串t。

StrConcat(s, t)

初始条件: 串s, t 已存在。

操作结果: 将串t联结到串s后形成新串存放到s中。

StrLength(t)

初始条件: 字符串t已存在。

操作结果: 返回串t中的元素个数, 称为串长。

SubString(s, pos, len, sub)

初始条件: 串s, 已存在, $1 \leq \text{pos} \leq \text{StrLength}(s)$ 且

$0 \leq \text{len} \leq \text{StrLength}(s) - \text{pos} + 1$ 。

操作结果: 用sub返回串s的第pos个字符起长度为len的子串。

.....

} ADT String



串的基本操作:



- (1) 创建串 **StringAssign**(s,string_constant)
- (2) 判断串是否为空 **StringEmpty**(s)
- (3) 计算串长度 **Length**(s)
- (4) 串连接 **Concat**(s1,s2)
- (5) 求子串 **SubStr**(s1,s2,start,len)
- (6) 子串的定位 **Index**(s1,s2)
- (7) 子串的插入和删除

2.4.2 串的存储设计



- **顺序串**:
 - 用数组来存储串中的字符序列。

- **非压缩形式**

c	h	i	n	e	s	e	
---	---	---	---	---	---	---	--

- **压缩形式**

c	e						
h	s						
i	e						
n							



如何表示串的长度？

- 方法一：用一个变量来表示串的实际长度，同一般线性表
- 方法二：在串尾存储一个不会在串中出现的特殊字符作为串的终止符，表示串的结尾。

0	1	2	3	4	5	6			Max
w	o	r	k	e	r	空	-----	闲	

0	1	2	3	4	5	6			Max
w	o	r	k	e	r	\0	空	-----	闲

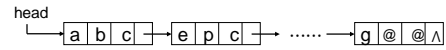
链接串：

用链接存储结构来存储串。

串的链式存储结构和线性表的串的链式存储结构类似，采用单链表来存储串，结点的构成是：

- data域：存放字符，data域可存放的字符个数称为结点的大小；
- next域：存放指向下一结点的指针。

若每个结点仅存放一个字符，则结点的指针域就非常多，造成系统空间浪费，为节省存储空间，考虑串结构的特殊性，使每个结点存放若干个字符，这种结构称为块链结构。



块大小为3串的块链式存储结构示意图



2.4.3 串的模式匹配算法 (Pattern Matching)

子串定位运算又称为模式匹配 (Pattern Matching) 或串匹配 (String Matching)，此运算的应用在非常广泛。例如，在文本编辑程序中，我们经常要查找某一特定单词在文本中出现的位置。显然，解此问题的有效算法能极大地提高文本编辑程序的响应性能。

在串匹配中，一般将主串称为目标串，子串称之为模式串。设 S 为目标串， T 为模式串，且不妨设：

$$S = "s_1s_2 \dots s_n" \quad T = "t_1 \dots t_m"$$

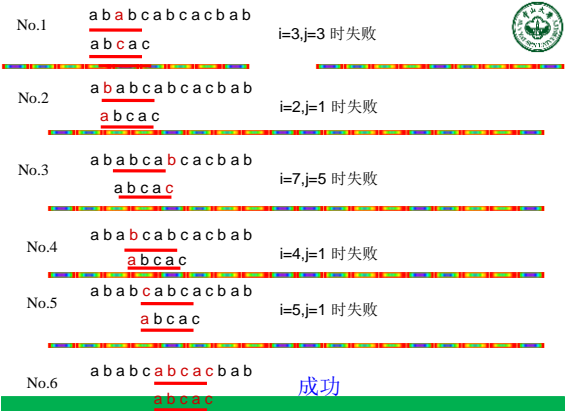
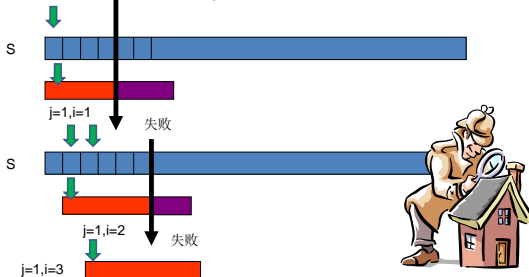
朴素模式匹配算法 (Brute-Force 算法)：枚举法

从主串 S 的第一个字符开始和模式 T 的第一个字符进行比较，若相等，则继续比较两者的后续字符；否则，从主串 S 的第二个字符开始和模式 T 的第一个字符进行比较，重复上述过程，直到 T 中的字符全部比较完毕，则说明本趟匹配成功；或 S 中字符全部比较完，则说明匹配失败。



其算法段为：

```
for (i=1; i<=n-m+1; i++) {
    if (S[i..i+m-1] == T[1..m])
        return i;
```



```

int index-1(ssstring s,ssstring t)
{
    int i,j,k;
    int n=s.length;
    int m=t.length;
    for(i=1;i<=n-m+1;i++)
    {
        j=1;k=i;
        while(j<=m && s.ch[k]==t.ch[j])
        {
            k++;j++;
        }
        if(j>m) return i;
    }
    return -1;
}

```



```

int index-2(ssstring s,ssstring t)
{
    int i=1,j=1;
    int n=s.length;
    int m=t.length;
    while (i<=n)&&(j<=m)
    {
        if ( s.ch[i]==t.ch[j])
        {
            i++;j++;
        }
        else
        {
            i=i-j+2;j=1;
        }
        if(j>m) return i-j+1;
        else return -1;
    }
}

```



时间复杂性

最好情况下算法的平均时间复杂性 $O(n+m)$ 。

最坏情况下的平均时间复杂性为 $O(n*m)$ 。

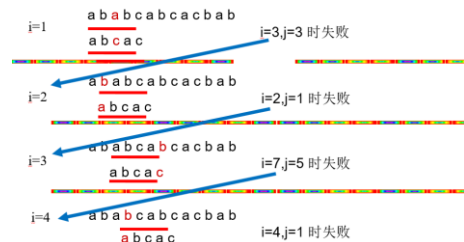


KMP 算法：改进的模式匹配算法



为什么BF算法时间性能低？

- 在每趟匹配不成功时存在大量回溯，没有利用已经部分匹配的结果。



KMP 算法：改进的模式匹配算法



为什么BF算法时间性能低？

- 在每趟匹配不成功时存在大量回溯，没有利用已经部分匹配的结果。

如何在匹配不成功时主串不回溯？

- 主串不回溯，模式就需要向右滑动一段距离。

如何确定模式的滑动距离？

- 利用已经得到的“部分匹配”的结果
- 将模式向右“滑动”尽可能远的一段距离 (next[j]) 后，继续进行比较

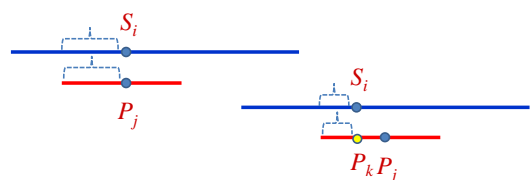
出发点：利用前面匹配的结果，进行无回溯匹配

如何在匹配不成功时主串不回溯？

- 主串不回溯，模式就需要向右滑动一段距离。

如何确定模式的滑动距离？

- 利用已经得到的“部分匹配”的结果
- 将模式向右“滑动”尽可能远的一段距离 (next[j]) 后，继续进行比较



KMP 算法



• 思考:

- 假定: 主串为 $S_1S_2\cdots S_n$
- 模式串为 $P_1P_2\cdots P_m$
- 当主串中的第 i 个字符和模式串中的第 j 个字符出现不匹配, 主串中的第 i 个字符应该和模式串中的哪个字符匹配(无回溯, 主串指针不回溯, 模式串右移动)?

• 进一步思考



- 假定主串中第 i 个字符与模式串第 j 个字符相比较失败, 则应有 $S_i \neq P_j$

$$\begin{array}{ccccccc}
 S_{i-j+1} & S_{i-j+2} & \cdots & S_{i-k+1} & S_{i-k+2} & \cdots & S_{i-1} & S_i \\
 P_1 & P_2 & \cdots & P_{j-k+1} & P_{j-k+2} & \cdots & P_{j-1} & P_j
 \end{array}$$

$k < j$

$$\begin{array}{ccccccc}
 S_{i-k+1} & S_{i-k+2} & \cdots & S_{i-1} & S_i \\
 P_1 & P_2 & \cdots & P_{k-1} & P_k
 \end{array}$$

成立

S_i 与 P_k 进行比较

$$\begin{array}{ccccccc}
 S_{i-j+1} & S_{i-j+2} & \cdots & S_{i-k+1} & S_{i-k+2} & \cdots & S_{i-1} & S_i \\
 P_1 & P_2 & \cdots & P_{j-k+1} & P_{j-k+2} & \cdots & P_{j-1} & P_j \\
 & & & P_1 & P_2 & \cdots & P_{k-1} & P_k
 \end{array}$$

- 而根据已有的匹配, 有

$$P_{j-k+1}P_{j-k+2}\cdots P_{j-1} = S_{i-k+1}S_{i-k+2}\cdots S_{i-1}$$

- 因此

$$P_{j-k+1}P_{j-k+2}\cdots P_{j-1} = P_1P_2\cdots P_{k-1}$$

- 因此 k 值只和 P 以及 j 有关, 定义为 $\text{Next}[j]$

意义: 当 P_j 比较失败时, 右移动模式串, 让 P_k 与当前 S_i 元素进行比较

• $\text{Next}[j]$ 的定义

$$\text{Next}[j] = \begin{cases} 0, & j=1 \text{ 时} \\ \text{Max}\{k \mid 1 < k < j \text{ and } p_1p_2\cdots p_{k-1} = p_{j-k+1}\cdots p_{j-1}\} & \\ 1, \text{其它情况} & \end{cases}$$

Next 数组的实质是找模式串中的最长相同的前缀和后缀。
 $p_1p_2\cdots p_{k-1} = p_{j-k+1}\cdots p_{j-1}$

j	1	2	3	4	5	6	7	8
	a	b	a	a	b	c	a	c
Next[j]	0	1	1	2	2	3	1	2

j	1	2	3	4	5
	a	b	c	a	c
Next[j]	0	1	1	1	2

No.1 a b a b c a b c a c b a b
 a b c a c i=3, j=3 时失败

No.2 a b a b c a b c a c b a b
 a b c a c i=7, j=5 时失败

No.3 a b a b c a b c a c b a b
 a b c a c 成功

```
int index(sstring s, sstring t)
```

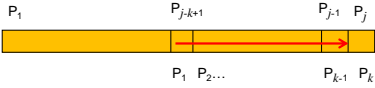
```
{
    int i=1, j=1;
    int n=s.length;
    int m=t.length;
    while (i<=n) && (j<=m)
    {
        if (s.ch[i]==t.ch[j]) (i=0) || (s.ch[i]==t.ch[j])
        {
            i++; j++;
        }
        else
        {
            i=i-j+2; j=1;
            j=Next[j];
        }
    }
    if (j>m) return i-j+1;
    else return -1;
}
```

计算Next数组的方法

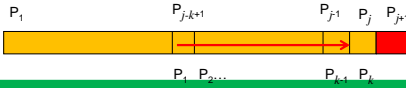
1) Next[1]=0;

2) 设 Next[j]=k; 则意味着

$$P_{j-k+1}P_{j-k+2}\dots P_{j-1}=P_1P_2\dots P_{k-1}$$

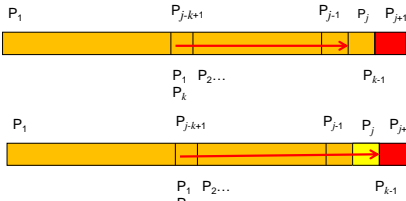


3) 求Next[j+1]



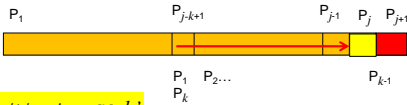
3) 求Next[j+1]

当 $P_k = P_j$ 时 $P_{j-k+1}P_{j-k+2}\dots P_{j-1} = P_1P_2\dots P_{k-1}$

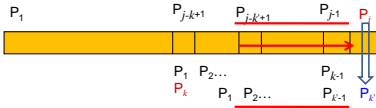


$$\text{Next}[j+1] = \text{Next}[j] + 1;$$

当 $P_k \neq P_j$ 时 $P_{j-k+1}P_{j-k+2}\dots P_{j-1} \neq P_1P_2\dots P_{k-1}$



存在一个next[k]=k'



当 $P_j = P_{k'}$, $\text{Next}[j+1] = \text{Next}[k'] + 1;$

同理, 若 $P_j = P_{k'}$, 则将模式串继续向右滑动至模式串的第next[k'] 个字符与 P_j 对齐, ..., 依次类推, 直到 P_j 和模式串中某个字符匹配成功或者不存在 k' ($1 \leq k' \leq j$), 则 $\text{Next}[j+1] = 1$.

Next数组的无回溯匹配计算

```
void Makenext(String p,int *pNext)
{
    int i,k; i=1;k=0;Next[1] = 0;
    while(i<Length(p))
    {
        if (k==0) || ( p[i]==p[k])
        { j=j+1; k=k+1;    Next[j]=k;  }
        else k=Next[k];
    }
}
```

Next数组特殊情况

j	1	2	3	4	5	6	7	8	9	10	11	12	13
a	a	a	a	a	a	a	b	b	b	c	a	a	
0	1	2	3	4	5	6	7	1	1	1	1	2	

改进后的Next数组:

j	1	2	3	4	5	6	7	8	9	10	11	12	13
a	a	a	a	a	a	a	b	b	b	c	a	a	
0	0	0	0	0	0	0	7	1	1	1	0	0	