

# KMP (**Knuth** Morris Pratt) 算法

《计算机程序设计艺术 第1卷 基本算法》

《计算机程序设计艺术 第2卷 半数值算法》

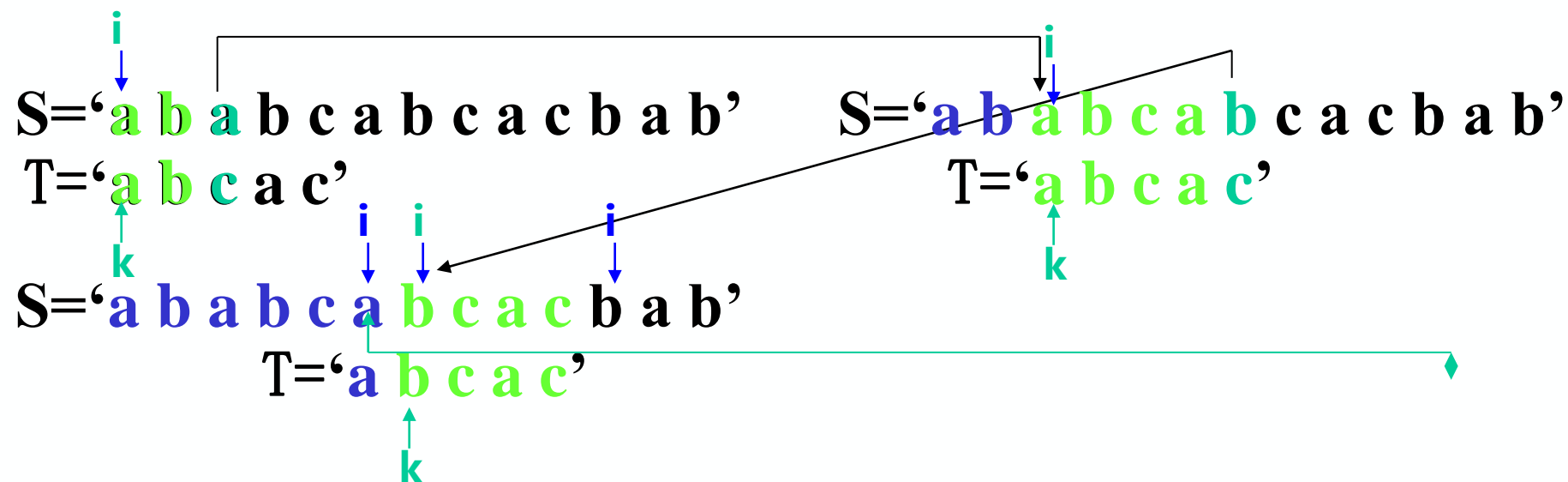
《计算机程序设计艺术 第3卷 排序与查找》

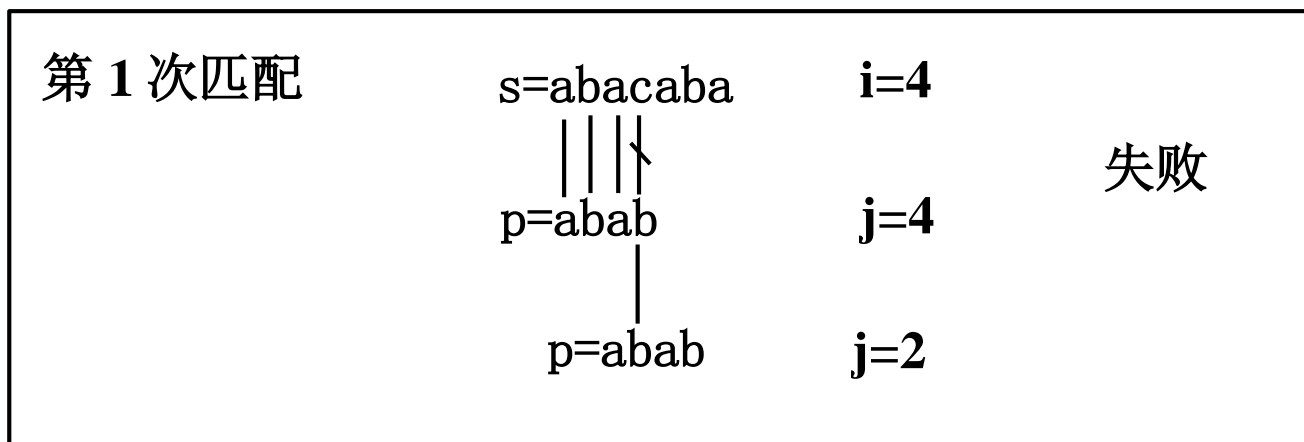
<http://www-cs-faculty.stanford.edu/~knuth/>



# KMP算法设计思想

利用已经**部分匹配**的结果而加快模式串的滑动速度？  
且主串S的指针**i不必回溯**！可提速到 **$O(n+m)$** ！





因 $p_1 \neq p_2$ ,  $s_2 = p_2$ , 必有 $s_2 \neq p_1$ , 又因 $p_1 = p_3$ ,  $s_3 = p_3$ ,  
所以必有 $s_3 = p_1$ 。因此, 第二次匹配可直接从  
 $i=4$ ,  $j=2$ 开始。

**改进：**每趟匹配过程中出现字符比较不等时，不回溯主指针 $i$ ，利用已得到的“部分匹配”结果将模式向右滑动尽可能远的一段距离，继续进行比较。

$$\begin{array}{ccccccc}
s_1 & s_2 & s_3 & \cdots & s_{i-j+1} & s_{i-j+2} & \cdots s_{i-2} & s_{i-1} & s_i & s_{i+1} \\
& & & & \parallel & \parallel & & \parallel & \parallel & \neq \\
& & & & p_1 & p_2 & \cdots & p_{j-2} & p_{j-1} & p_j & p_{j+1} \\
& & & & & & \parallel & & \parallel & & \\
& & & & & & p_1 & \cdots & p_{k-1} & p_k & p_{k+1}
\end{array}$$

① “ $p_1 p_2 \cdots p_{k-1}$ ” = “ $s_{i-k+1} s_{i-k+2} \cdots s_{i-1}$ ”

② “ $p_{j-k+1} p_{j-k+2} \cdots p_{j-1}$ ” = “ $s_{i-k+1} s_{i-k+2} \cdots s_{i-1}$ ” (部分匹配)

③ “ $p_1 p_2 \cdots p_{k-1}$ ” = “ $p_{j-k+1} p_{j-k+2} \cdots p_{j-1}$ ” (真子串)

为此,定义 $\text{next}[j]$ 函数, 表明当模式中第 $j$ 个字符与主串中相应字符“失配”时, 在模式中需重新和主串中该字符进行比较的字符的位置。

$$\text{next}[j]=\begin{cases} \max\{ k|1 < k < j, \text{且 } "p_1 \cdots p_{k-1}" = "p_{j-k+1} \cdots p_{j-1}" \} & \text{当此集合非空时} \\ 0 & \text{当 } j=1 \text{ 时} \\ 1 & \text{其他情况} \end{cases}$$

```
int Index_KMP (SString S,SString T, int pos)
{
    i= pos,j =1;
    while (i<S[0] && j<T[0]) {
        if (j==0 || S[i]==T[j]) {  i++;j++; }
        else
            j=next[j];          /*i不变,j后退*/
    }
    if (j>T[0]) return i-T[0]; /*匹配成功*/
    else return 0;             /*返回不匹配标志*/
}
```

## ● 如何求next函数值

1.  $\text{next}[1] = 0$ ; 表明主串从下一字符 $s_{i+1}$ 起和模式串重新开始匹配。  $i = i+1; j = 1$ ;

2. 设 $\text{next}[j] = k$ , 则 $\text{next}[j+1] = ?$

① 若 $p_k = p_j$ , 则有“ $p_1 \cdots p_{k-1} p_k$ ”=“ $p_{j-k+1} \cdots p_{j-1} p_j$ ”, 如果在

$j+1$ 发生不匹配, 说明 $\text{next}[j+1] = k+1 = \text{next}[j]+1$ 。

② 若 $p_k \neq p_j$ , 可把求next值问题看成是一个模式匹配问题, 整个模式串既是主串, 又是子串。



$$\begin{array}{ccccccc}
 p_1 & p_2 & \cdots & p_{j-k+1} & \cdots & p_{j-1} & p_j & p_{j+1} & \text{next}[j]=k \\
 & & & \parallel & & \parallel & \neq & & \\
 & & & p_1 & \cdots & p_{k-1} & p_k & p_{k+1} & \text{next}[k]=k' \\
 & & & p_1 & \cdots & p_k & p_{k'+1} & & \text{next}[k']=k'' \\
 & & & p_1 & \cdots & p_{k''} & p_{k''+1} & & \text{next}[k'']=k'''
 \end{array}$$

- 若  $p_{k'}=p_j$ , 则有 “ $p_1 \cdots p_{k'}$ ” = “ $p_{j-k'+1} \cdots p_j$ ”,  
 $\text{next}[j+1]=k'+1=\text{next}[k]+1=\text{next}[\text{next}[j]]+1$ .
- 若  $p_{k''}=p_j$ , 则有 “ $p_1 \cdots p_{k''}$ ” = “ $p_{j-k''+1} \cdots p_j$ ”,  
 $\text{next}[j+1]=k''+1=\text{next}[k'] + 1 = \text{next}[\text{next}[k]] + 1$ .
- $\text{next}[j+1]=1$ .

<b>j</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>10</b>	<b>11</b>	<b>12</b>	<b>13</b>	<b>14</b>	<b>15</b>	<b>16</b>	<b>17</b>
<b>模式串</b>	<b>a</b>	<b>b</b>	<b>c</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>b</b>	<b>c</b>	<b>a</b>	<b>b</b>	<b>c</b>	<b>a</b>	<b>a</b>	<b>b</b>	<b>d</b>	<b>a</b>	<b>b</b>
<b>next[j]</b>	<b>0</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>2</b>	<b>2</b>	<b>3</b>	<b>1</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>1</b>	<b>2</b>

- **void get\_next(SString T, int &next[])**  
**{**  
    **i= 1; next[1] = 0; j = 0;**  
    **while( i<T.length){**  
        **if(j==0 || T[i] == T[j]){**  
            **++i; ++j;**  
            **next[i] = j;**  
        **}**  
        **else**  
            **j = next[j];**  
    **}**  
**}**

## ● KMP算法的时间复杂度

设主串  $s$  的长度为  $n$ , 模式串  $t$  长度为  $m$ , 在KMP算法中求  $next$  数组的时间复杂度为  $O(m)$ , 在后面的匹配中因主串  $s$  的下标不减即不回溯, 比较次数可记为  $n$ , 所以KMP算法总的时间复杂度为  $O(n+m)$ 。