



北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

数据结构与算法设计

课程内容



课程内容：数据结构部分

概述

线性表

栈与队列

数组与广义表

串

树

图

查找

内部排序

外部排序



课程内容：算法设计部分

概述

分治

动态规划

贪心

回溯

.....

.....

.....

.....

计算模型

可计算理论

计算复杂性



串的基本概念

串的基本操作

KMP算法

BM算法

Sunday算法

一、什么是串

串是有零个或多个字符组成的有限序列。一般记为：

$$s = 'a_1 a_2 \dots a_n' \quad \text{其中 } n \geq 0$$

串的概念

一、什么是串

$s = 'a_1 a_2 \dots a_n'$,其中 $n \geq 0$

s 称为串的**名**, $'\dots'$ 之间的字符序列称为串的**值**

a 可以是字母, 数字或其他字符

字符的数目称为串的长度

零个字符的串, 称为空串, 长度为0

串的概念

一、什么是串

$s = 'a_1 a_2 \dots a_n'$,其中 $n \geq 0$

串中任意连续个字符构成的子序列称为**子串**

字符在串中的序列中的序号, 称为**字符的位置**

子串第一个字符在主串中的位置, 称为**子串的位置**

a= 'Bei' b= 'Jing'

c= 'BeiJing' d= 'Bei Jing'

a,b都是c和d的子串。

a的长度是3， b的长度是4， c的长度是7， d的长度是8

a在c， d中的位置都是0， b在c中的位置是3， d中的位置是4

当且仅当两个串的值相等，我们说这两个串是相等的

换句话说，当两个串的字符序列和它们在字符序列中的位置都相等，那么这两个串是相等的。

空格串，由一个或多个空格组成的串，叫做空格串。

$s = \quad ' \quad '$

空串是任何串的子串，通常记为 \emptyset

一、特殊的线性表

串与线性表极度相似

线性表的操作常常是单个元素为处理对象的；

串的操作常常是以串或子串为处理对象的；

串含有结束标记；

1) 用常量初始化串StrAssign(&T,chars)

前提：chars是一个字符串常量

功能：创建一个值等于chars的串T

2) 字符串复制Strcpy(&T, S)

前提：S存在

功能：由串S复制得串T

3) 判别S是否为空串StrEmpty(S)

前提：S必须存在

功能：若S为空串，返回True，否则返回false

串的基本操作

4) 字符串比较 **StrCompare (T,S)**

前提: T, S存在

功能: 若 $S > T$, 返回值 > 0 , 若S与T相等, 则返回值 $= 0$,
若 $S < T$, 返回值 < 0

5) 求字符串长度 **StrLength(S)**

功能: 返回串S的字符个数, 即S的长度

6) 清空串 **ClearStr(&S)**

前提: S必须存在

功能: 将S清空为空串

串的基本操作

7) 字符串连接 **ConcatStr (&T, S1, S2)**

前提: S1, S2存在

功能: 连接S1和S2组成新的串T

8) 求子串 **SubStr(&Sub, S, pos, len)**

前提: 串S存在, $0 \leq \text{pos} \leq \text{StrLength}(S)$ 且

$0 \leq \text{len} \leq \text{StrLength}(S) - \text{pos}$

功能: 返回串S的第pos个字符起, 长度为len的子串Sub

9) 返回子串第一次出现的位置 `getIndex(S, T, pos)`

前提：S, T必须存在, T是非空串, $0 \leq \text{pos} \leq \text{StrLength}(S)$

功能：若S中含有和串T相同的子串, 则返回它在串S中的第pos个字符之后第一次出现的位置; 否则返回-1

10) 子串替换 `replace (&S, T, V)`

前提：S, T和V都存在, T是非空串

功能：用V替换串S中出现的所有与T相等的不重叠子串

串的基本操作

11) 插入子串 `StrInsert (&S, pos, V)`

前提: S和T存在, $0 \leq \text{pos} \leq \text{StrLength}(S)$

功能: 在串S的第pos个位置之前插入串T

12) 删除子串 `StrDelete (&S, pos, len)`

前提: S存在, $0 \leq \text{pos} \leq \text{StrLength}(S) - \text{len}$

功能: 从S中删除第pos个字符起长度为len的子串

13) 释放串的存储空间 (`&S`)

前提: S存在

功能: 串S被销毁

二、串的物理实现方式

- 顺序结构（数组，C语言）
- 链表方式（略）

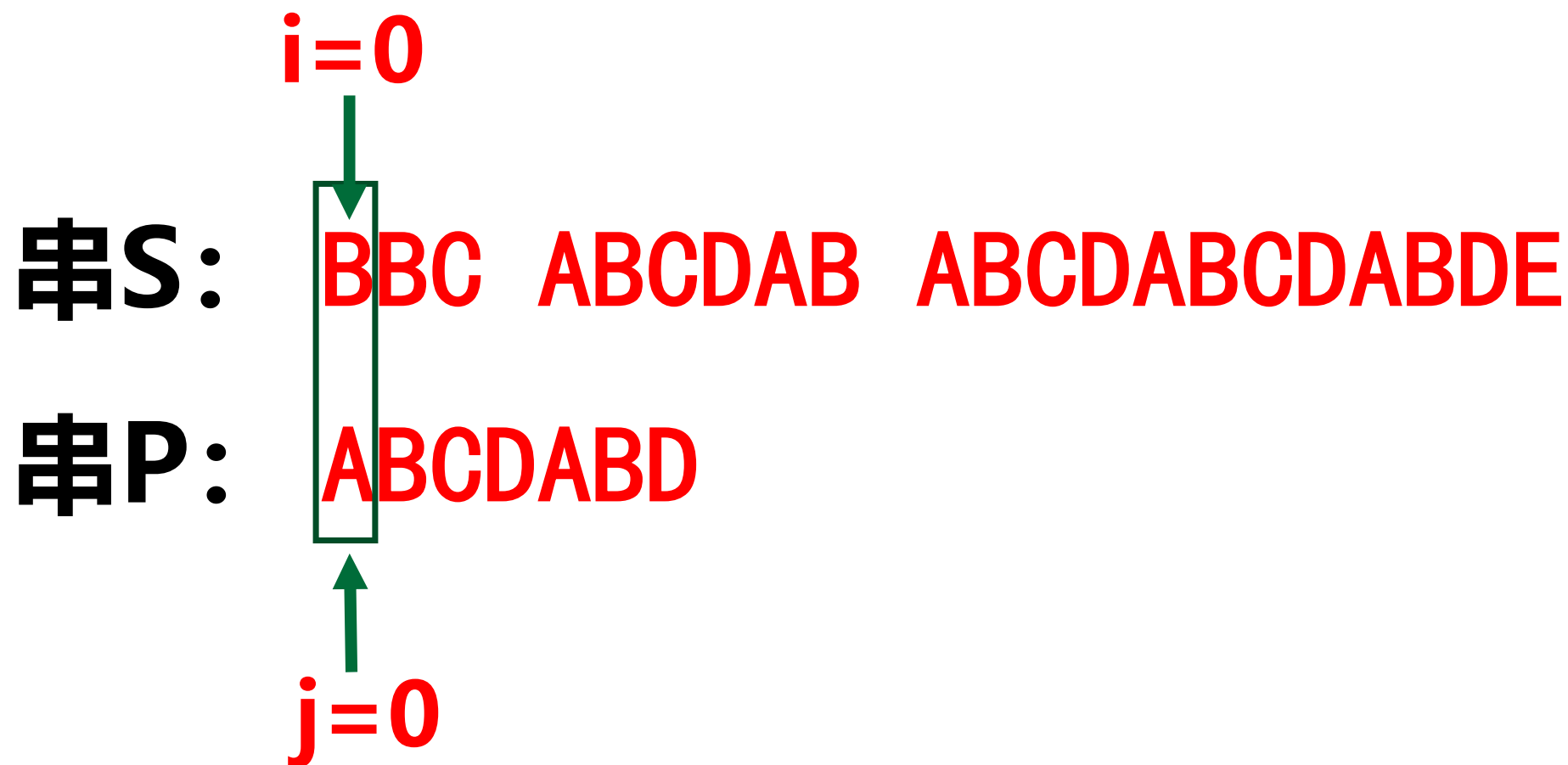
串的模式匹配算法

串的模式匹配算法：查找子串定位的Index(S,P,pos)_暴力法

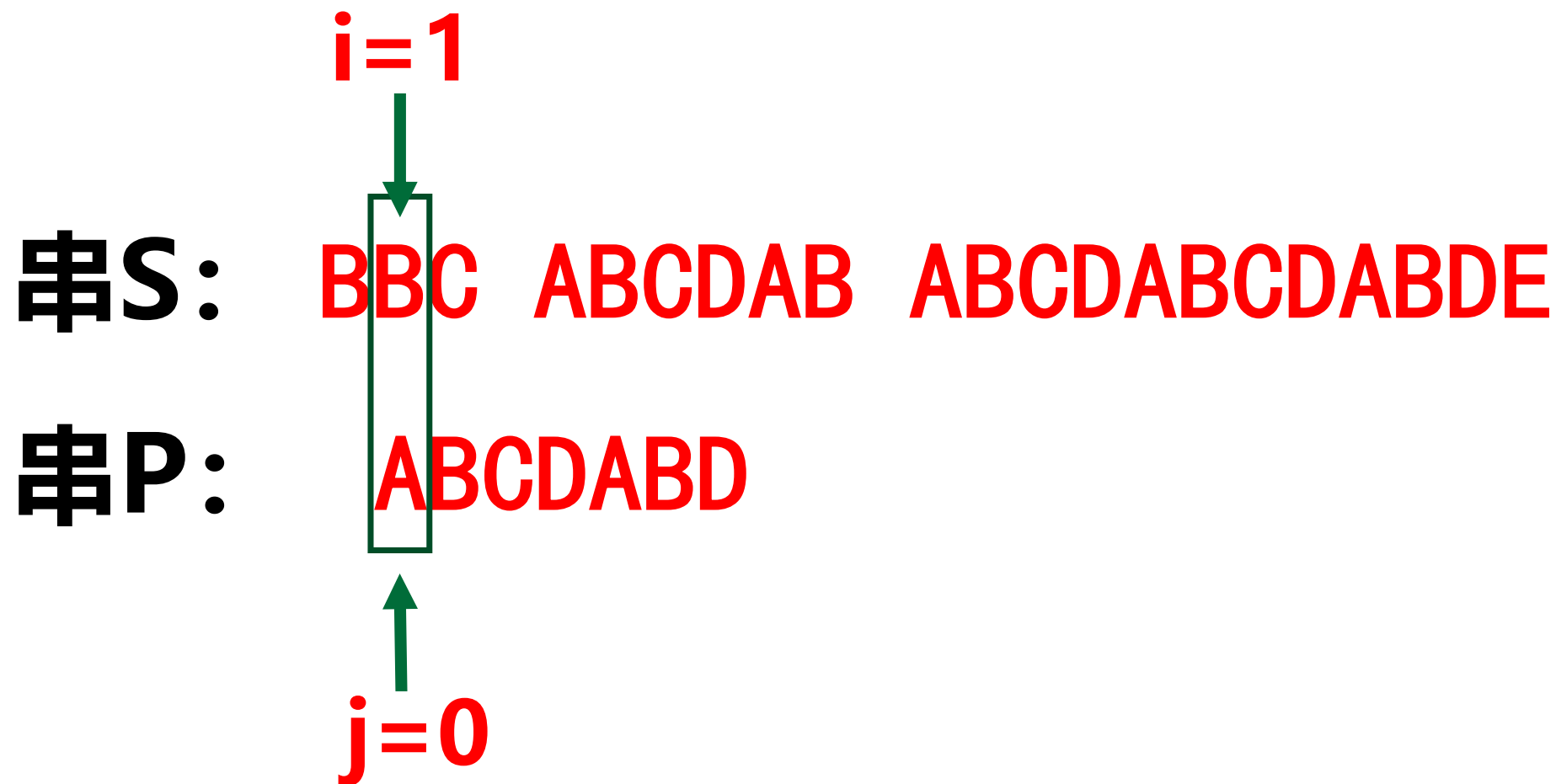
```
int i=pos , j=0;
while(i<strlen(S)&& j<strlen(P))
{
    if ( S[i] == P[j] ){ i++; j++; }
    else { i=i-j+1; j=0; }
}
if (i>=strlen(S)) printf( "%d" ,-1); //匹配不成功
else printf("%d",i-j); //匹配成功
```

最好的匹配成功的情况是m,
最好的匹配不成功的情况是n,
最坏的情况是 $n*m$

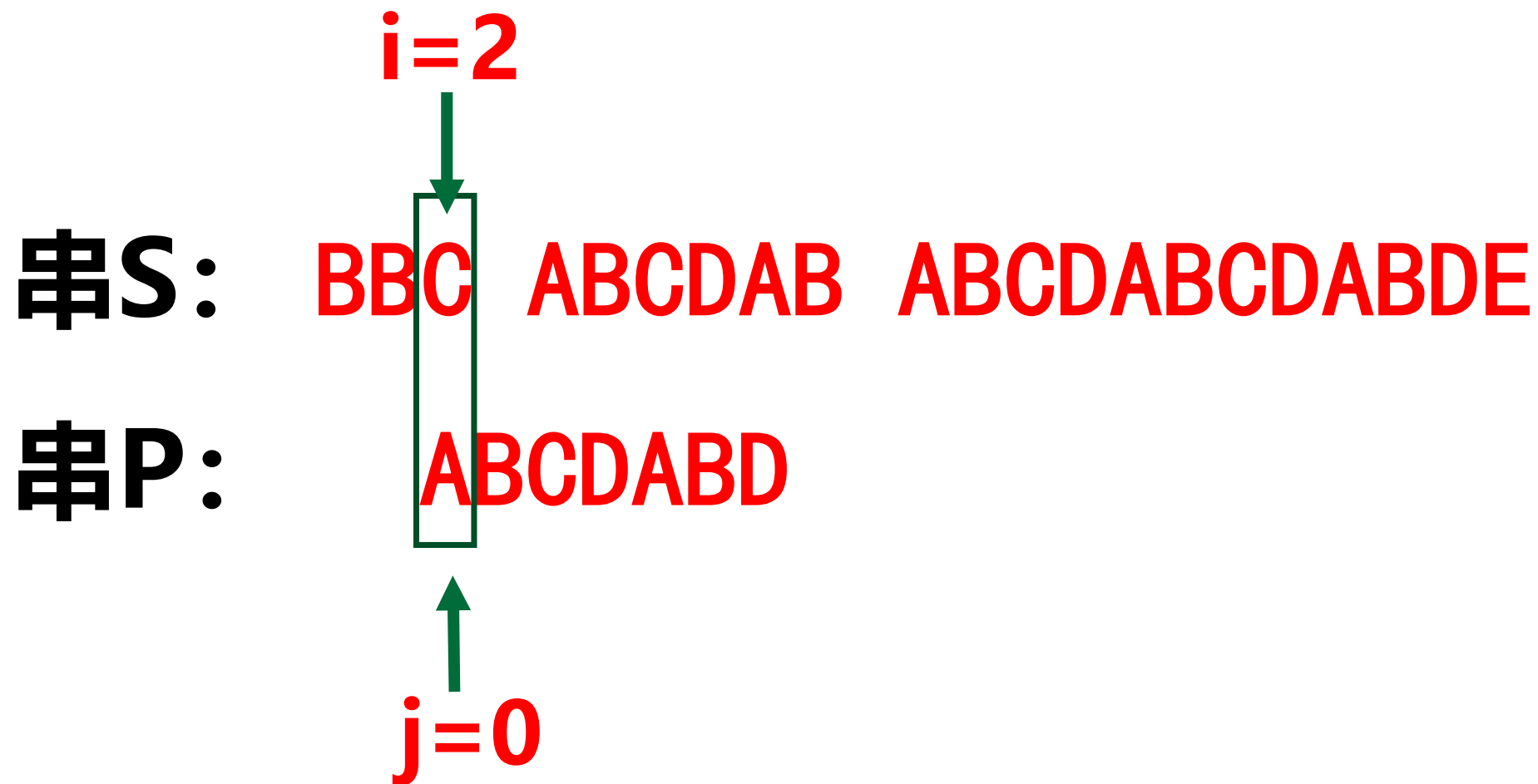
串的模式匹配算法

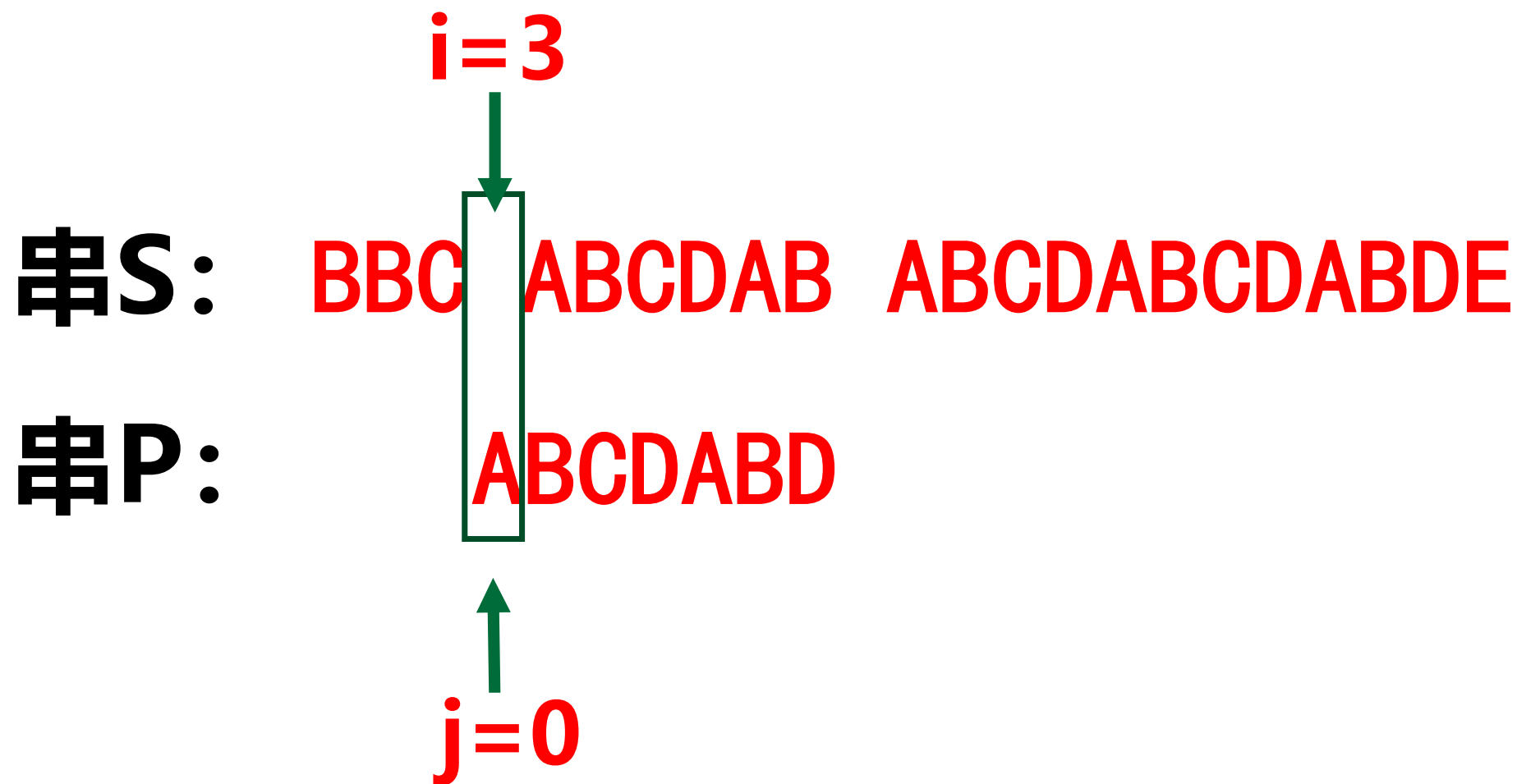


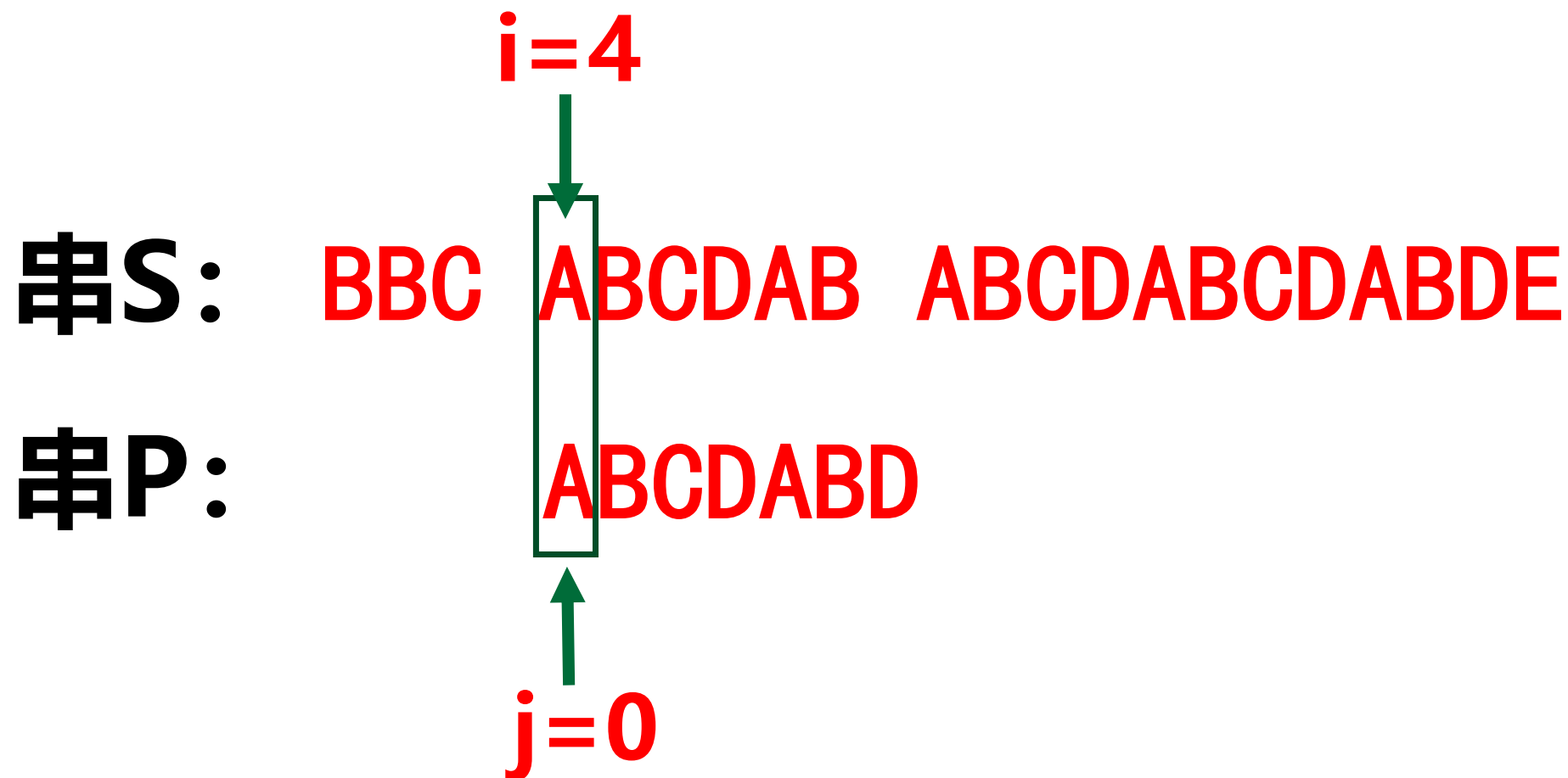
串的模式匹配算法

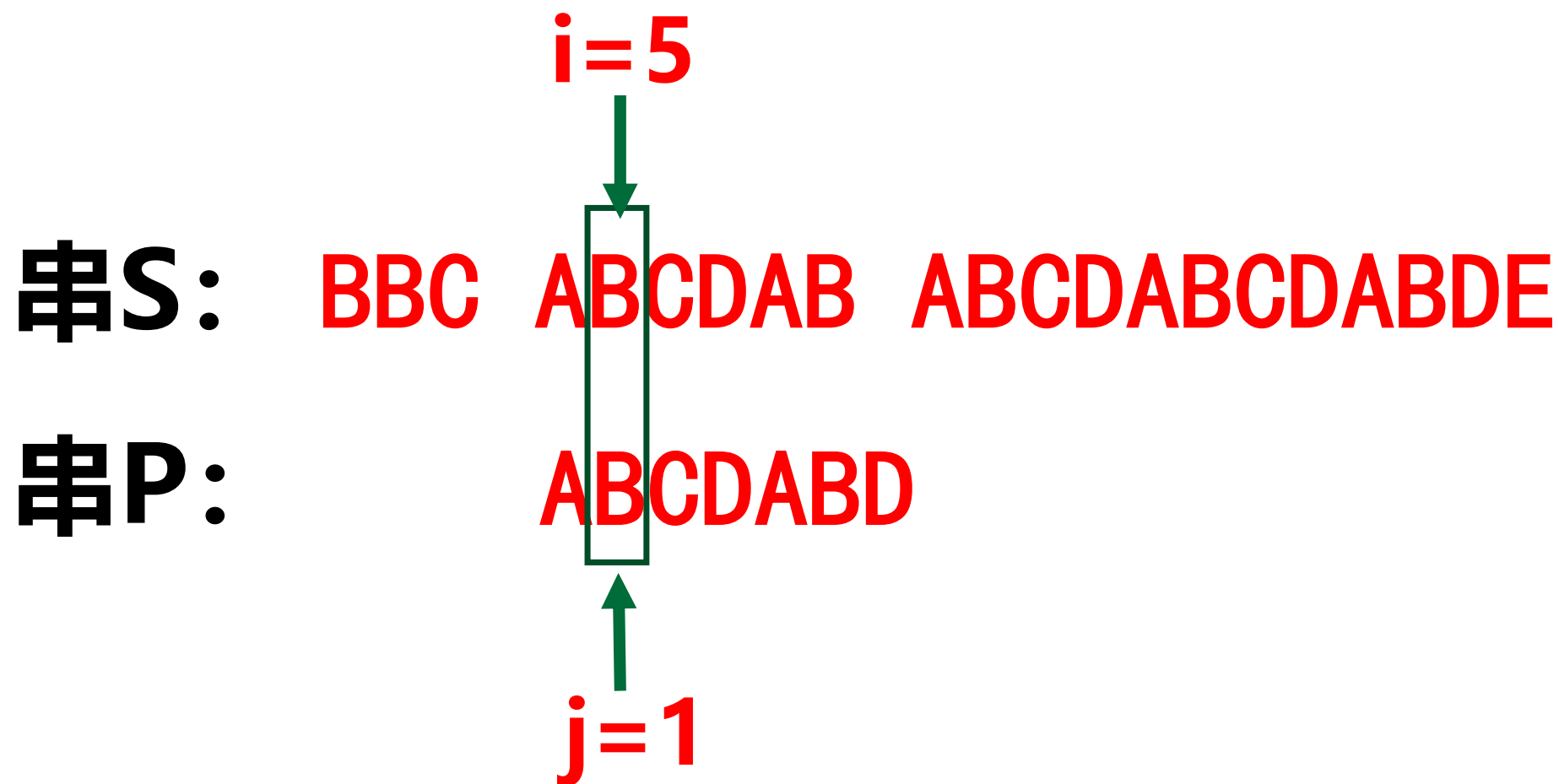


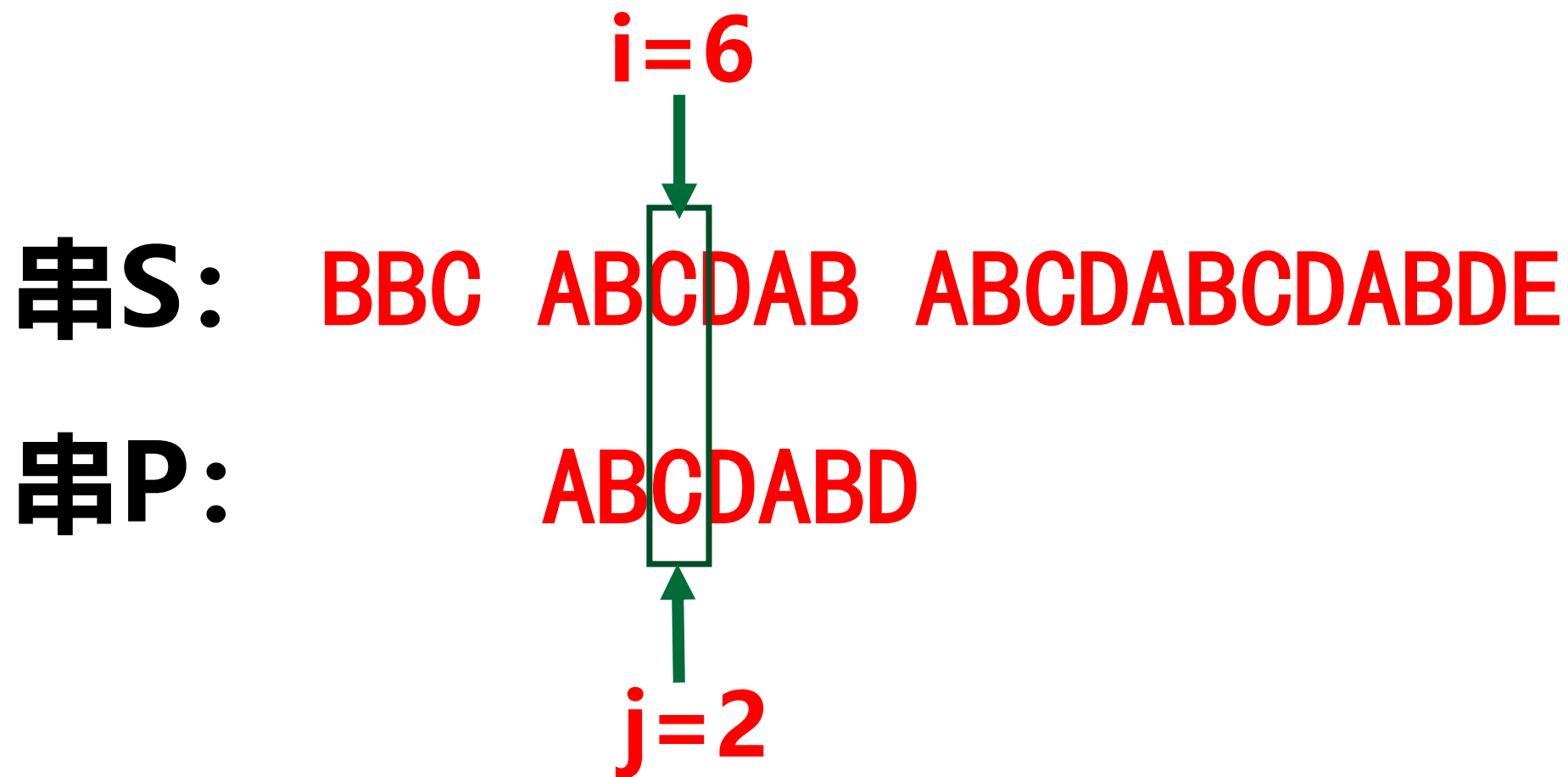
串的模式匹配算法

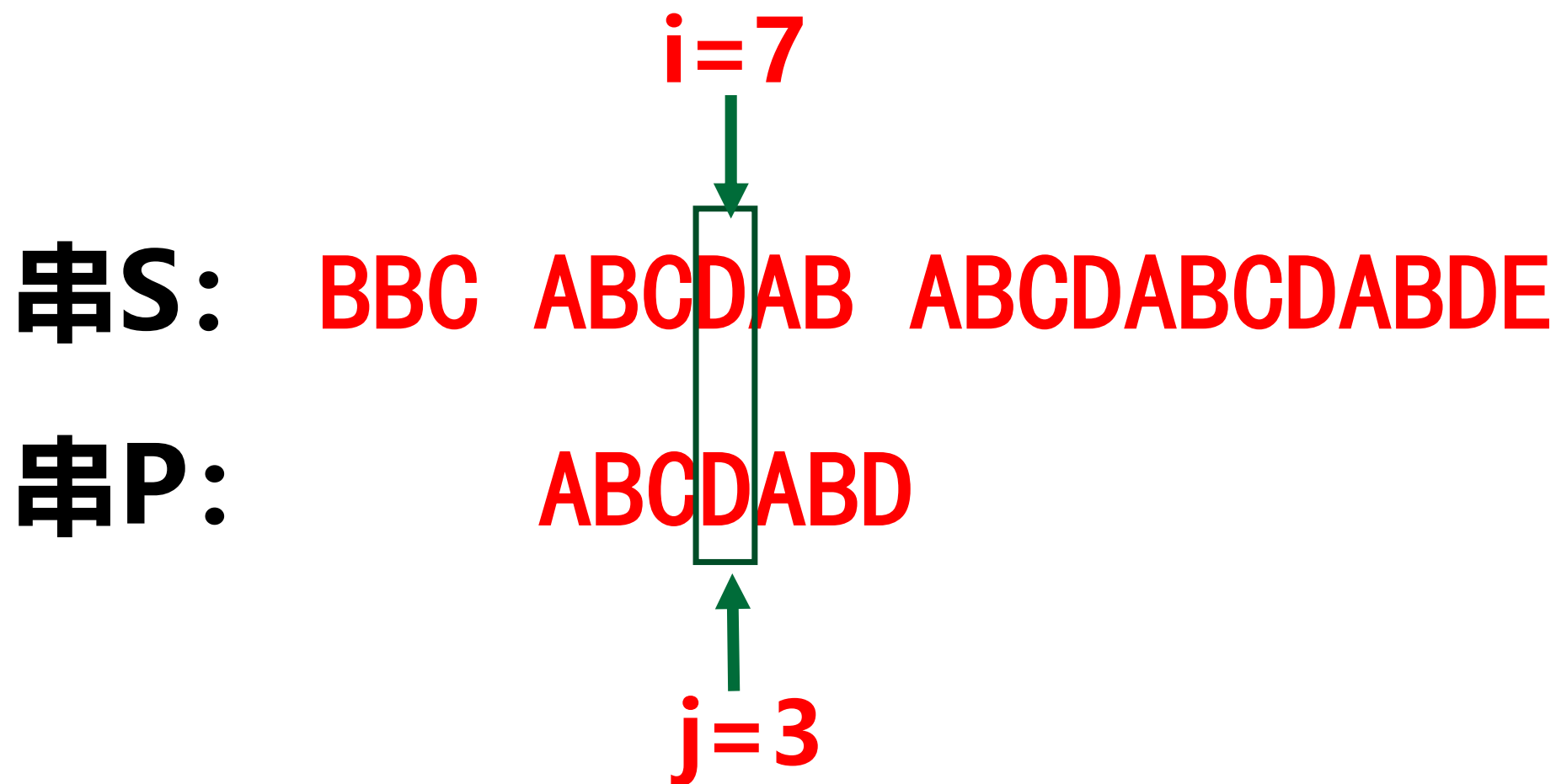


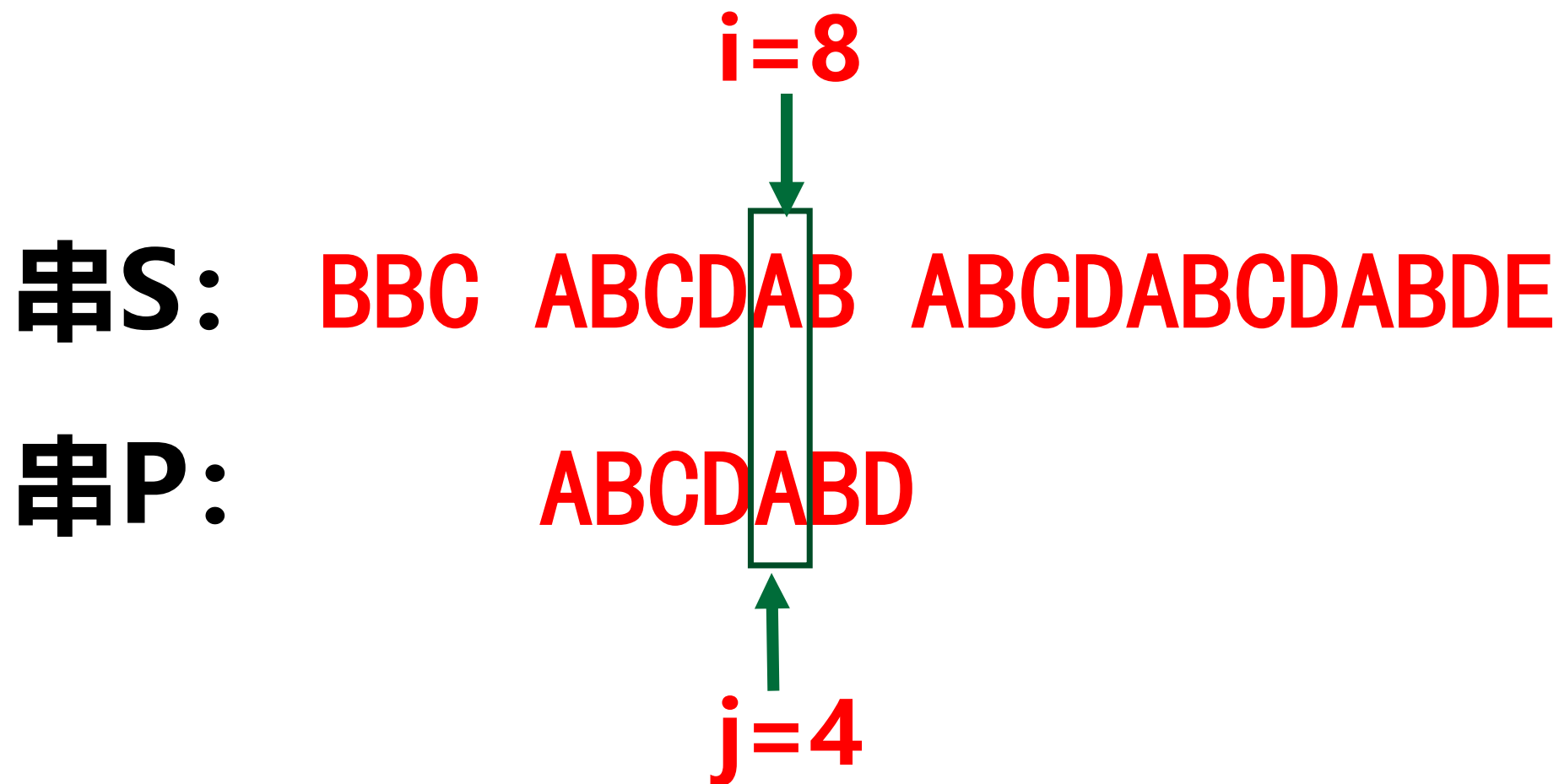


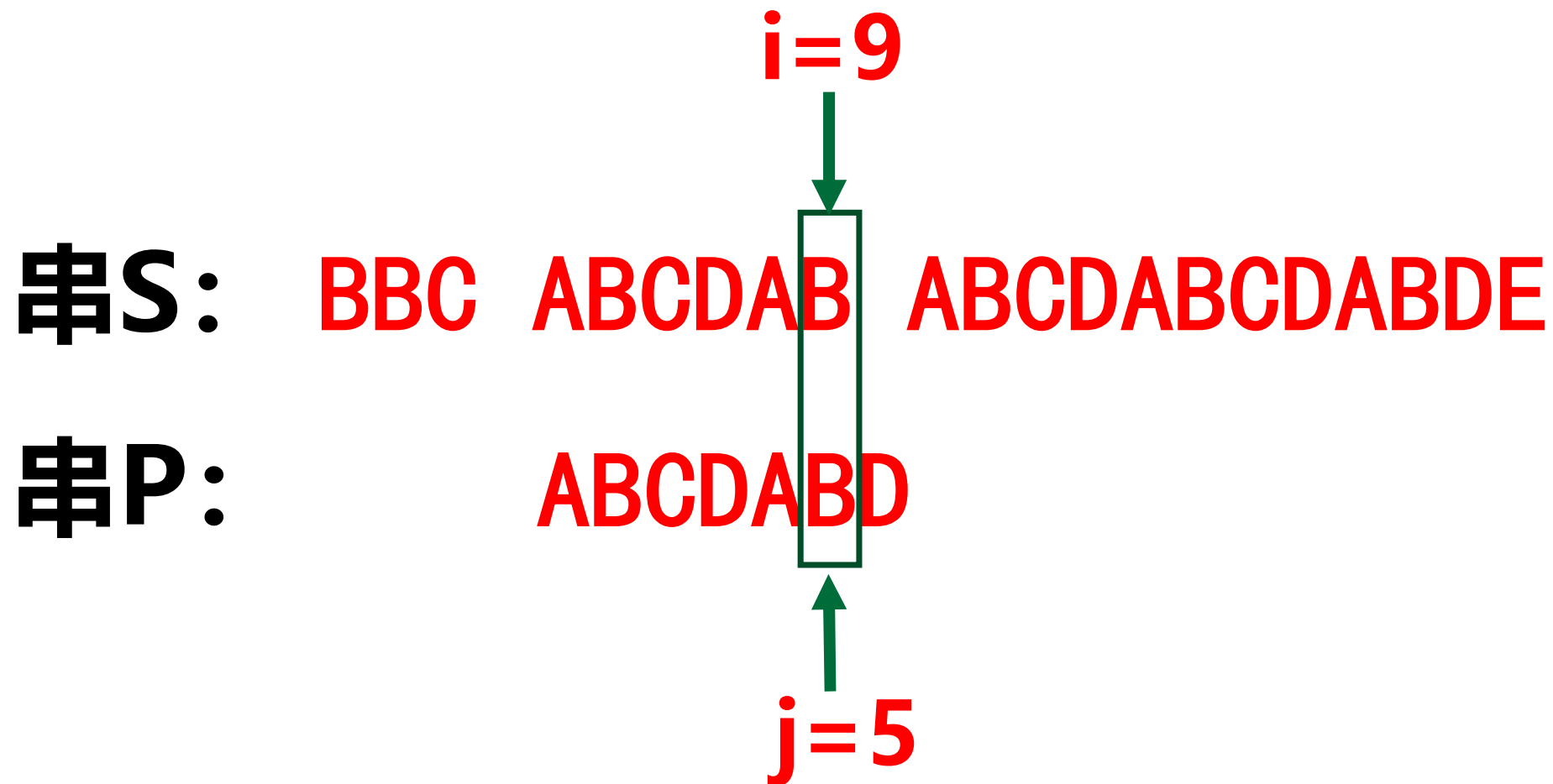


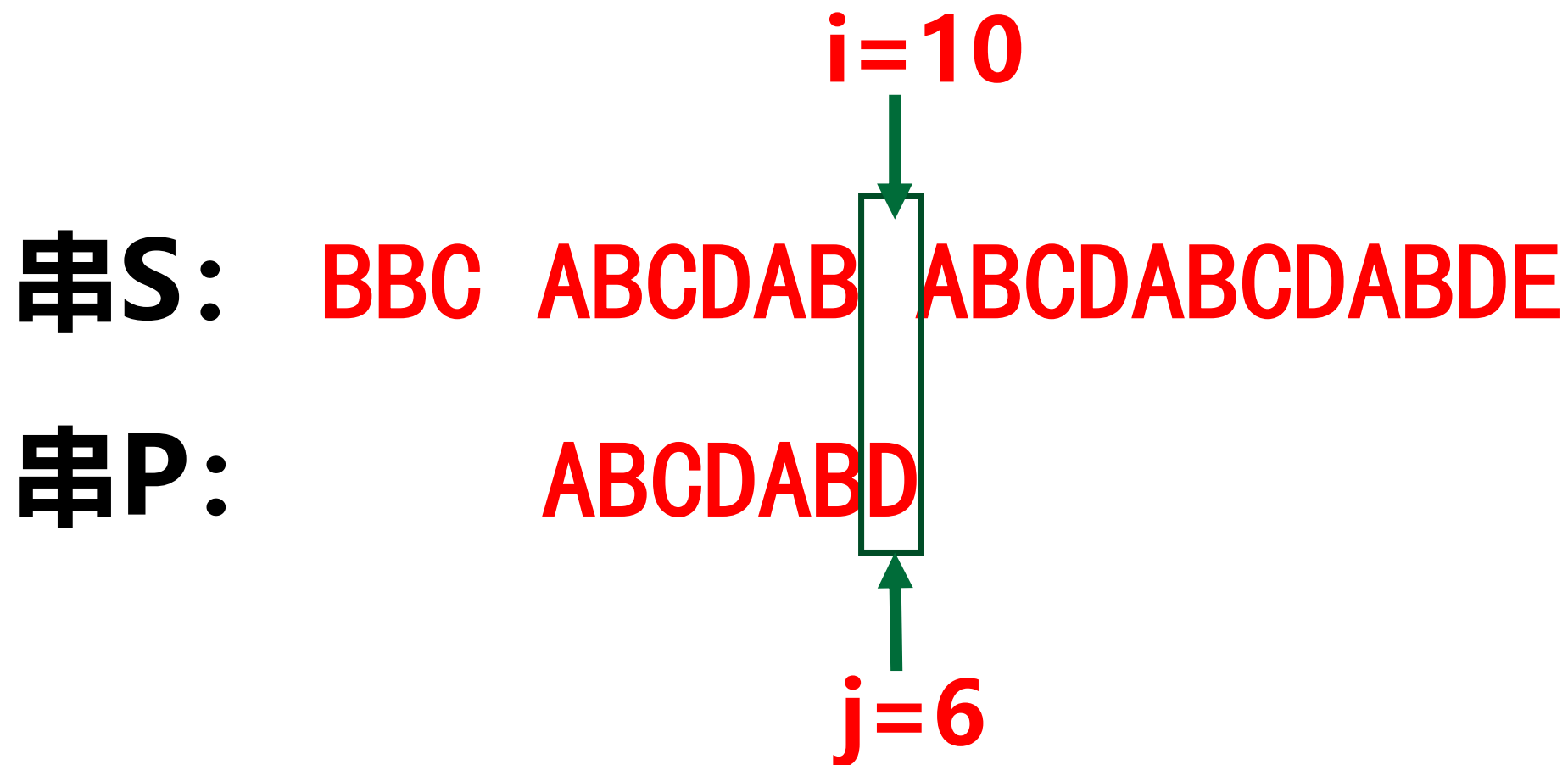


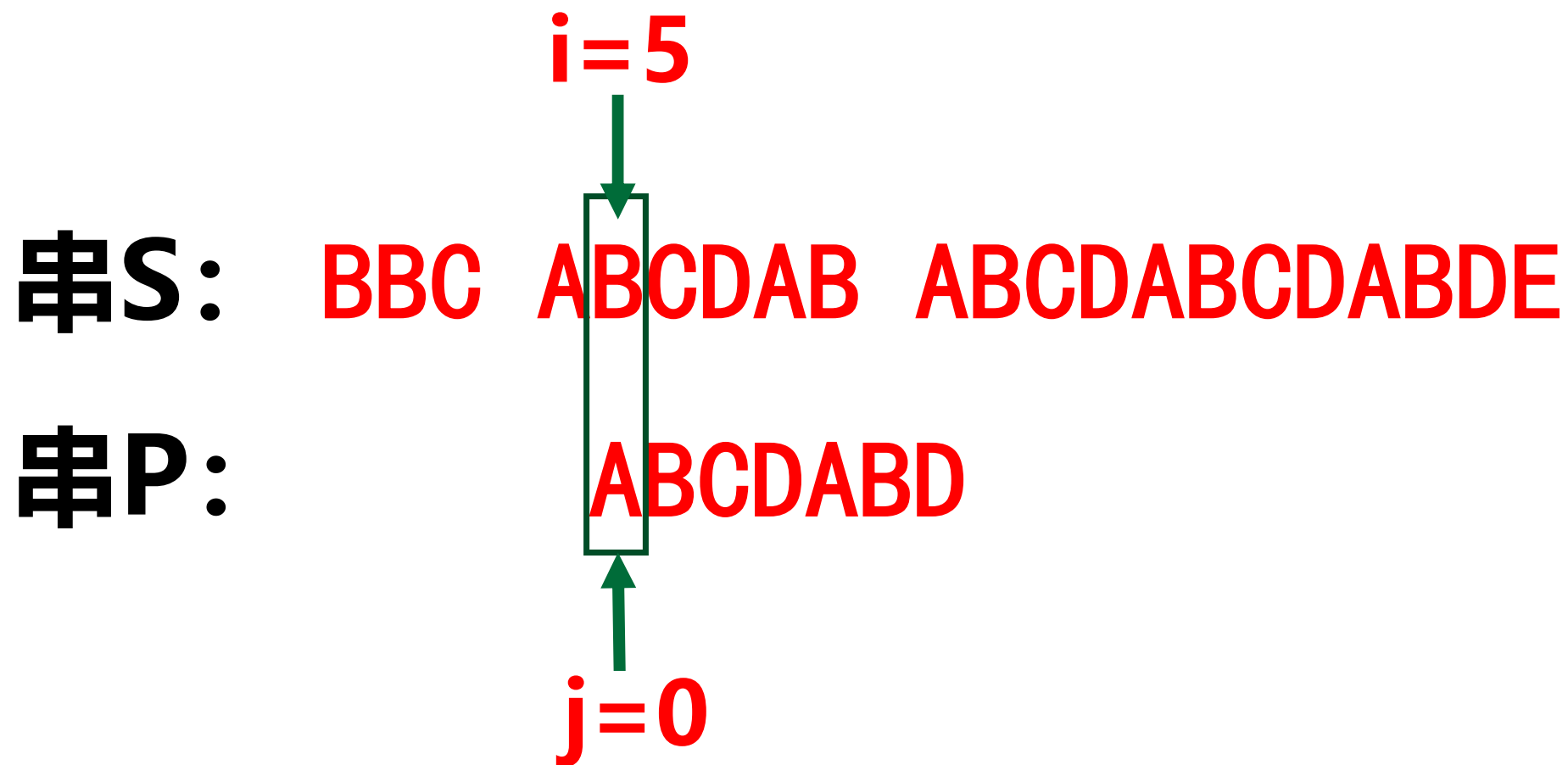












串的模式匹配算法：KMP算法

算法由D.E.**K**nuth、J.H.**M**orris和V.R.**P**ratt,共同发现

利用之前已经部分匹配这个有效信息，保持**i**不回溯，通过修改**j**的位置，让模式串尽量地移动到有效的位置。

串的模式匹配算法：KMP算法

假设现在文本串S匹配到 i 位置，模式串P匹配到 j 位置：

如果当前匹配失败：

在 i 保持不变的情况下， j 应该从哪里开始进行匹配？

串的模式匹配算法：KMP算法

假设现在文本串S匹配到 i 位置，模式串P匹配到 j 位置：

如果当前匹配失败：

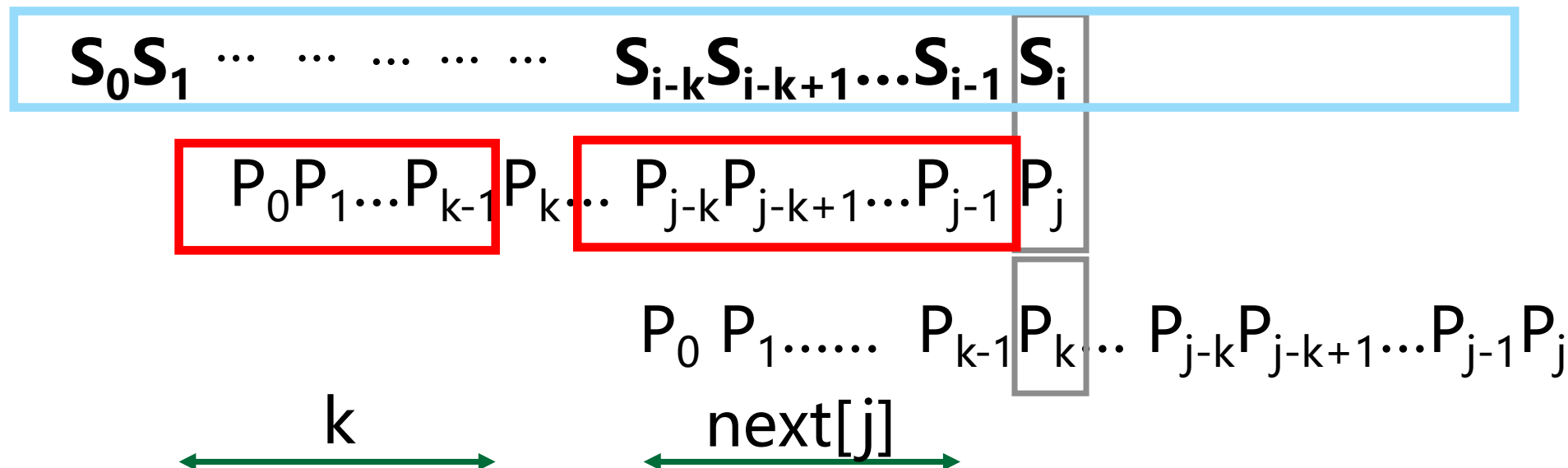
在 i 保持不变的情况下， j 应该从哪里开始进行匹配？

假定 i 位置的字符，应该与P串中的第 k 个字符做匹配，那么：

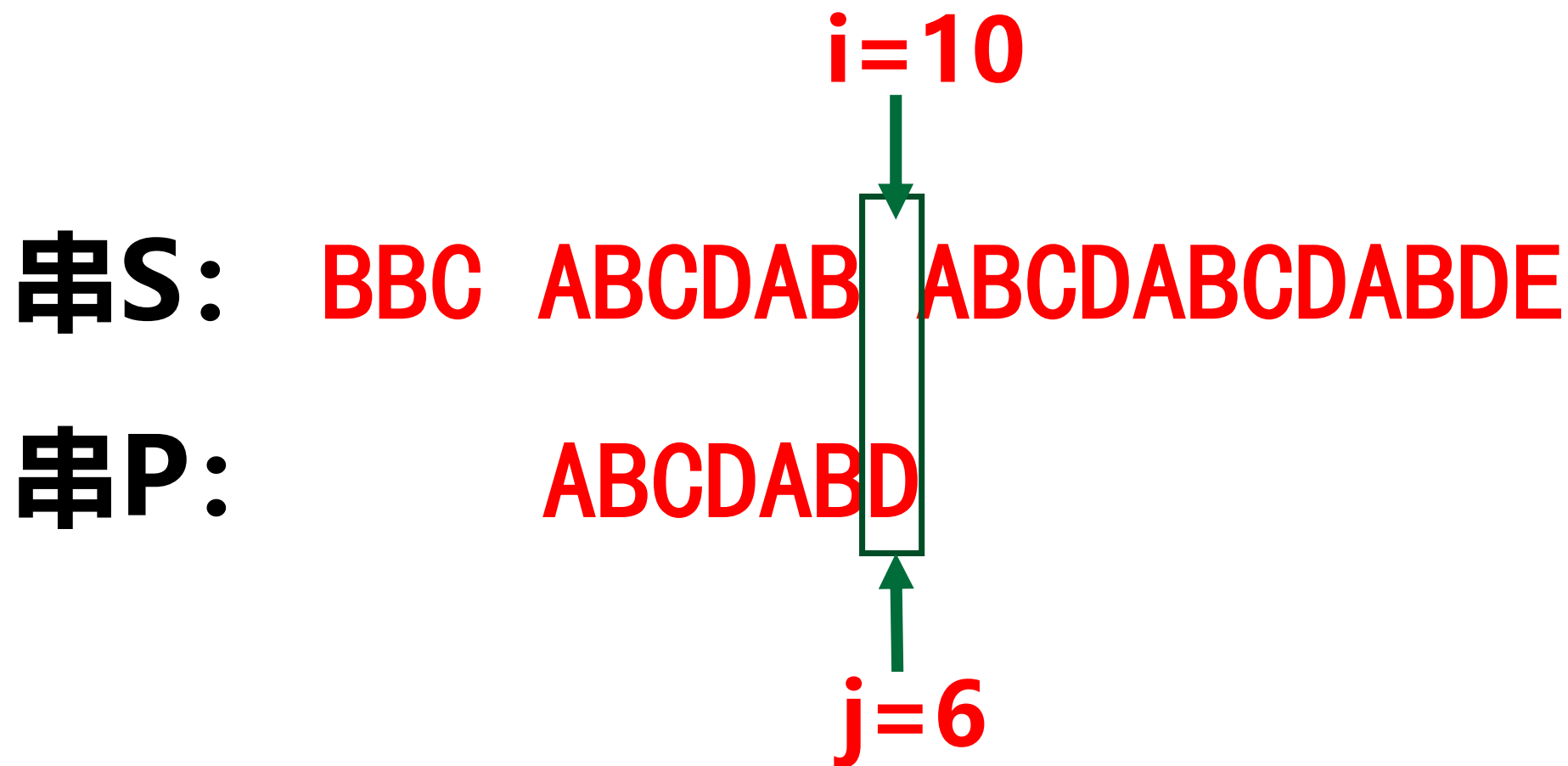
串的模式匹配算法：KMP算法

$S_0 S_1 \dots \dots \dots S_{i-k} S_{i-k+1} \dots S_{i-1} S_i$
$P_0 P_1 \dots P_{k-1} P_k \dots P_{j-k} P_{j-k+1} \dots P_{j-1} P_j$

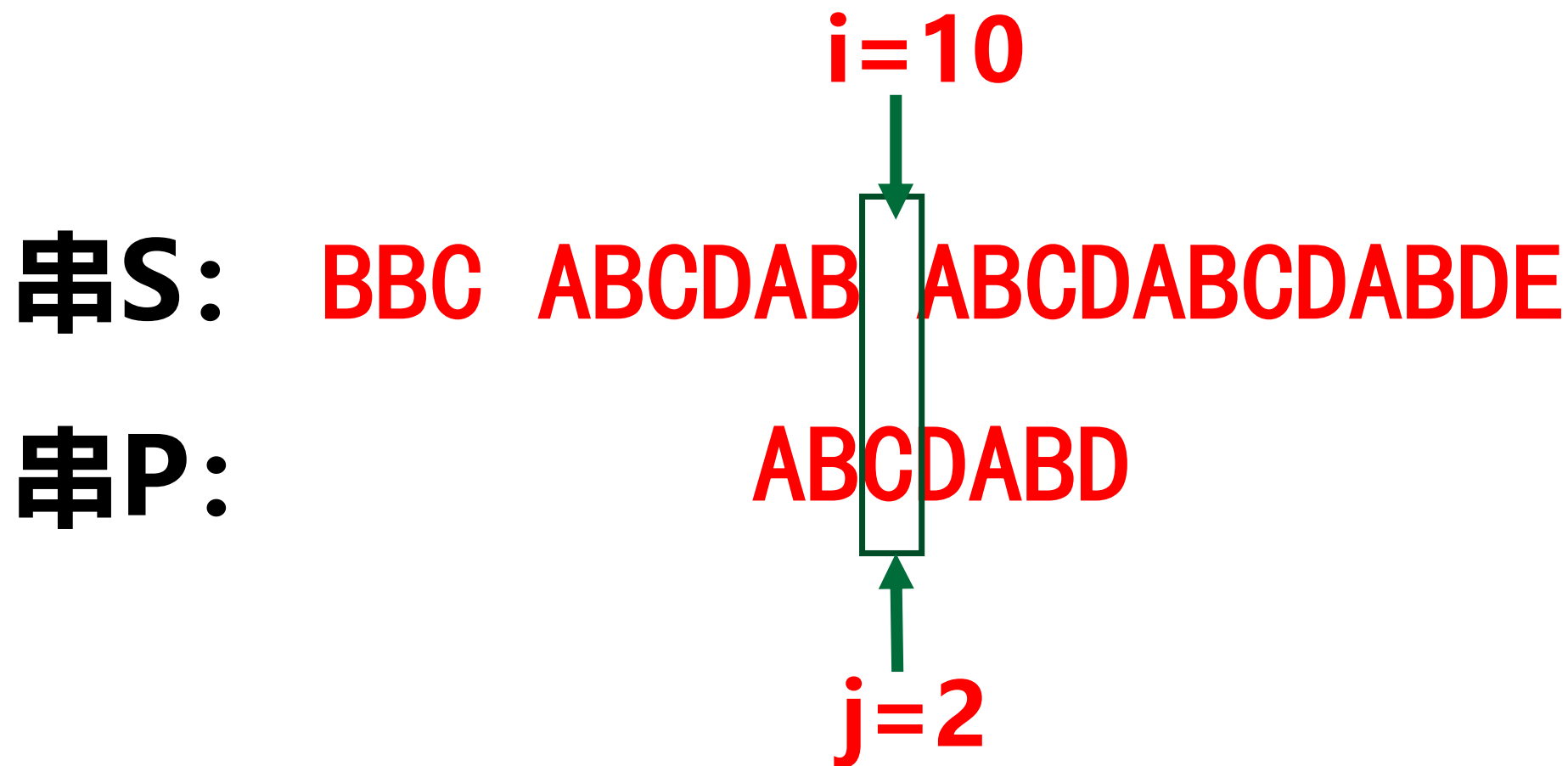
串的模式匹配算法：KMP算法



串的模式匹配算法：KMP算法



串的模式匹配算法：KMP算法



串的模式匹配算法：KMP算法

定义Next数组，用于存储当i与j匹配失败的时候，下一步应该让i去匹配P[k]字符，那么如何求取next数组呢？

- (1) 求取最大前缀后缀公共元素长度；
- (2) 向右平移1位，左侧补-1；

串的模式匹配算法：KMP算法

模式串的各个子串	前缀	后缀	最大公共元素长度
A	空	空	0
AB	A	B	0
ABC	A,AB	C,BC	0
ABCD	A,AB,ABC	D,CD,BCD	0
ABCDA	A,AB,ABC,ABCD	A,DA,CDA,BCDA	1
ABCDAB	A,AB,ABC,ABCD,ABCDA	B,AB,DAB,CDAB,BCDAB	2
ABCDABD	A,AB,ABC,ABCD,ABCDA ABCDAB	D,BD,ABD,DABD,CDABD BCDABD	0

串的模式匹配算法：KMP算法

next 数组考虑的是除当前字符外的最长相同前缀后缀，所以通过第①步骤求得各个前缀后缀的公共元素的最大长度后，只要稍作变形即可：将第①步骤中求得的值整体右移一位，然后初值赋为-1，如下表格所示：

模式串	A	B	C	D	A	B	D
最大前缀后缀 公共元素长度	0	0	0	0	1	2	0
j	0	1	2	3	4	5	6
next数组的值： next[j]	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

定义Next数组，用于存储当i与j匹配失败的时候，下一步应该让i去匹配P[k]字符，关键是求得这个k：

问题变成了，求解k的过程：

(1) 给出k的定义求解：

$$\text{next}[j] = \begin{cases} -1 & \text{当 } j=0 \text{ 的时候} \\ \text{Max} \{ k \mid 0 < k < j \text{ 且 } "P_0P_1\dots P_{k-1}P_k" == "P_{j-k}P_{j-k+1}\dots P_{j-1}" \} & \\ 0 & \text{其他情况} \end{cases}$$

串的模式匹配算法：KMP算法

(3) 给出k的递归求解过程：

令： $\text{next}[0] = -1, \text{next}[1] = 0$;

设： $\text{next}[j] = k$ ，即表示当j字符“失配”时，应退到第k个元素，
因为：有 “ $P_0P_1\dots P_{k-1}$ ” == “ $P_{j-k}P_{j-k+1}\dots P_{j-1}$ ”

接下来，考察 $\text{next}[j+1]$ ：

若 $P_k = P_j$ ，则有 “ $P_0P_1\dots P_{k-1}P_k$ ” == “ $P_{j-k}P_{j-k+1}\dots P_{j-1}P_j$ ”，显然
也有 $\text{next}[j+1] = k+1$ ，

若 $P_k \neq P_j$ ，则有两种情况：

a) $\text{next}[j+1] = \text{next}[k] + 1$ ，（若 $P[\text{next}[k]] = P[j]$ ）

b) if $(k == -1)$ $\text{next}[j+1] = 0$ else $k = \text{next}[k]$ ，

串的模式匹配算法：KMP算法

ABCDABCE



j

$j == 0$:

$\text{next}[0] = -1$, 算法固定

串的模式匹配算法：KMP算法

ABCDABCE



j

$j=1$:

$\text{next}[1] = 0$, k值为0, 算法固定

串的模式匹配算法：KMP算法

ABCDABCE



$\text{next}[2] = 0$, k 的值为0

$P_k(\text{A}) \neq P_j(\text{C})$, 则有两种情况:

a) $\text{next}[j+1] = \text{next}[k] + 1$, (若 $P[\text{next}[k]] = P[j]$)

b) if $(k == -1)$ $\text{next}[j+1] = 0$ else $k = \text{next}[k]$

推算出 $\text{next}[3] = 0$

串的模式匹配算法：KMP算法

ABCDABCE



$\text{next}[3] = 0$, k 的值为0

$P_k(A) \neq P_j(D)$, 则有两种情况:

- a) $\text{next}[j+1] = \text{next}[k] + 1$, (若 $P[\text{next}[k]] = P[j]$)
- b) if $(k = -1)$ $\text{next}[j+1] = 0$ else $k = \text{next}[k]$,
推算出 $\text{next}[4] = 0$

串的模式匹配算法：KMP算法

ABCDABCE



$\text{next}[4] = 0$, k 的值为0

$P_k(A) == P_j(A)$,

若 $P_k == P_j$, 则有 " $P_0P_1\dots P_{k-1}P_k$ " == " $P_{j-k}P_{j-k+1}\dots P_{j-1}P_j$ ", 显然也有 $\text{next}[j+1] = k+1$,

则: $\text{next}[5] = 1$

串的模式匹配算法：KMP算法

ABCDABCE



$\text{next}[5] = 1$, k 的值为1

$P_k(\text{B}) = P_j(\text{B})$,

若 $P_k = P_j$, 则有 “ $P_0P_1\dots P_{k-1}P_k$ ” == “ $P_{j-k}P_{j-k+1}\dots P_{j-1}P_j$ ”, 显然也有 $\text{next}[j+1] = k+1$,

则: $\text{next}[6] = 2$

串的模式匹配算法：KMP算法

ABCDABCE



$\text{next}[6] = 2$, k 的值为2

$P_k(\text{C}) == P_j(\text{C})$,

若 $P_k = P_j$, 则有 " $P_0P_1\dots P_{k-1}P_k$ " == " $P_{j-k}P_{j-k+1}\dots P_{j-1}P_j$ ", 显然也有 $\text{next}[j+1] = k+1$,

则: $\text{next}[7] = 3$

串的模式匹配算法：KMP算法

ABCDABCE



$\text{next}[7] = 3$, k的值为3

串的模式匹配算法：KMP算法

ABCDABCE

模式串	A	B	C	D	A	B	C	E
j	0	1	2	3	4	5	6	7
next数组的 值：next[j]	-1	0	0	0	0	1	2	3

串的模式匹配算法：KMP算法

next 数组递推过程（初始值）

模式串	A	B	C	D	A	B	D
j	0	1	2	3	4	5	6
next数组的值： next[j]	-1	0					

根据 P_0, P_1 , 求next[2]的值:

P_j 是 'B', k 是0, P_k 是 'A', 所以next[2]为0

串的模式匹配算法：KMP算法

next 数组递推过程 ($\text{next}[2]=0$)

模式串	A	B	C	D	A	B	D
j	0	1	2	3	4	5	6
next数组的值: next[j]	-1	0	0				

根据 P_0, P_1, P_2 求 $\text{next}[3]$ 的值:

P_j 是 'C' , k 是0, P_k 是 'A' , 所以 $\text{next}[3]$ 为0

串的模式匹配算法：KMP算法

next 数组递推过程 (next[3]=0)

模式串	A	B	C	D	A	B	D
j	0	1	2	3	4	5	6
next数组的值: next[j]	-1	0	0	0			

根据 P_0, P_1, P_2, P_3 求next[4]的值:

P_j 是 'D' , k 是0, P_k 是 'A' , 所以next[4]为0

串的模式匹配算法：KMP算法

next 数组递推过程 (next[4]=0)

模式串	A	B	C	D	A	B	D
j	0	1	2	3	4	5	6
next数组的值: next[j]	-1	0	0	0	0		

根据 P_0, P_1, P_2, P_3, P_4 , 求next[5]的值:
 P_j 是 'A', k 是0, P_k 是 'A', 所以next[5]为next[4]+1

串的模式匹配算法：KMP算法

next 数组递推过程 (next[5]=1)

模式串	A	B	C	D	A	B	D
j	0	1	2	3	4	5	6
next数组的值: next[j]	-1	0	0	0	0	1	

根据 $P_0, P_1, P_2, P_3, P_4, P_5$ 求next[6]的值:
 P_j 是 'B' , k 是1, P_k 是 'B' , 所以next[6]为next[5]+1

串的模式匹配算法：KMP算法

next 数组递推过程 (next[6]=2)

模式串	A	B	C	D	A	B	D
j	0	1	2	3	4	5	6
next数组的值: next[j]	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0							

根据 P_0, P_1 ，求next[2]的值：

P_j 是 'A'， k 是0， P_k 是 'D'，所以next[2]为0

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0	0						

根据 P_0 , P_1 , P_2 求next[3]的值：

P_j 是 'B' , k 是0, P_k 是 'D' , 所以next[3]为0

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0	0	0					

根据 P_0, P_1, P_2, P_3 ，求next[4]的值：

P_j 是 'C'， k 是0， P_k 是 'D'，所以next[4]为0

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0	0	0	0				

根据 P_0, P_1, P_2, P_3, P_4 ，求next[5]的值：
 P_j 是 'D'， k 是0， P_k 是 'D'，所以next[5]为 $k+1$ ，
即 $0+1=1$

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0	0	0	0	1			

根据 $P_0, P_1, P_2, P_3, P_4, P_5$ ，求next[6]的值：
 P_j 是 'A'， k 是1， P_k 是 'A'，所以next[6]为 $k+1$ ，
即 $1+1=2$

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0	0	0	0	1	2		

根据 $P_0, P_1, P_2, P_3, P_4, P_5, P_6$ ，求next[7]的值：
 P_j 是 'B'， k 是2， P_k 是 'B'，所以next[7]为 $k+1$ ，
即 $2+1=3$

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0	0	0	0	1	2	3	

根据 $P_0, P_1, P_2, P_3, P_4, P_5, P_6, P_7$ ，求next[8]的值：
 P_j 是 'D'， k 是3， P_k 是 'C'，
所以next[8]为： $\text{next}[3] + 1 = 1$

串的模式匹配算法：KMP算法

next 数组递推过程示例3：

模式串	D	A	B	C	D	A	B	D	E
j	0	1	2	3	4	5	6	7	8
next数组的 值：next[j]	-1	0	0	0	0	1	2	3	1

串的模式匹配算法：KMP算法

next数组的求解过程

```
int j,k;
k=-1;j=0;
next[0]=-1;
while(j<S.length){
    if (k==-1 || (S[j]==S[k])) {
        j++;
        k++;
        next[j]=k;
    } else {
        k=next[k];
    }
}
```

串的模式匹配算法：KMP算法

(2) next数组的使用 (初始 $i=0$, $j=0$) :

在匹配过程中:

(0) 比较 S_i 和 P_j ;

(1) 若 $S_i = P_j$, 则 $i++$, $j++$, 继续 (0) 匹配 S_i 和 P_j ;

(2) 若 $S_i \neq P_j$, 则退回到 k ($\text{next}[j]$) 的位置, 即 j 赋值为 k ($k < j$), 再比较:

若 S_i 与 $P_{k(j)}$ 相等: 则 $i++$, $j++$;

若 S_i 与 $P_{k(j)}$ 不相等, 则继续退回到 $\text{next}[k=j]$ 的位置,

$\text{next}[k]$ 变相等于 $\text{next}[\text{next}[j]]$, 依次类推下去, 直到:

a:退回到 $\text{next}[\dots\text{next}[j]]$;

b:退回到0, 则开始比较 S_{i+1} 与 P_0

串的模式匹配算法：KMP算法

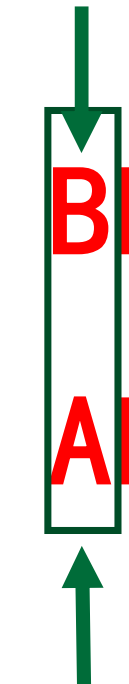
模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD

i=0

j=0



串的模式匹配算法：KMP算法

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

$i=1$



串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD



$j=0$

串的模式匹配算法：KMP算法

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

$i=2$



串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD



$j=0$

串的模式匹配算法：KMP算法

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

i=3



串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD



j=0

串的模式匹配算法：KMP算法

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

i=4

串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD

j=0

串的模式匹配算法：KMP算法

串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD

$i=10$

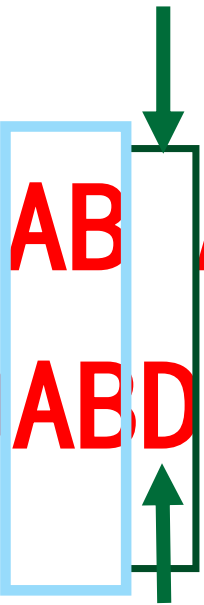
$j=6$

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD



$i=10$

$j=2$

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD

$i=10$

$j=0$

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

串S: BBC ABCDAB ABCDABCDABDE

串P: ABCDABD

$i=10$

$j=0$

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

串S: BBC ABCDAB **AB**CDABCDABDE

串P: **A**BCDABD

$i=11$

$j=0$

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

串S: BBC ABCDAB ABCDABCDABDE

串P:

ABCDABD

$i=17$

$j=6$

模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

串S: BBC ABCDAB ABCDABCDABDE

串P:

$i=17$
ABCDABD
 $j=2$

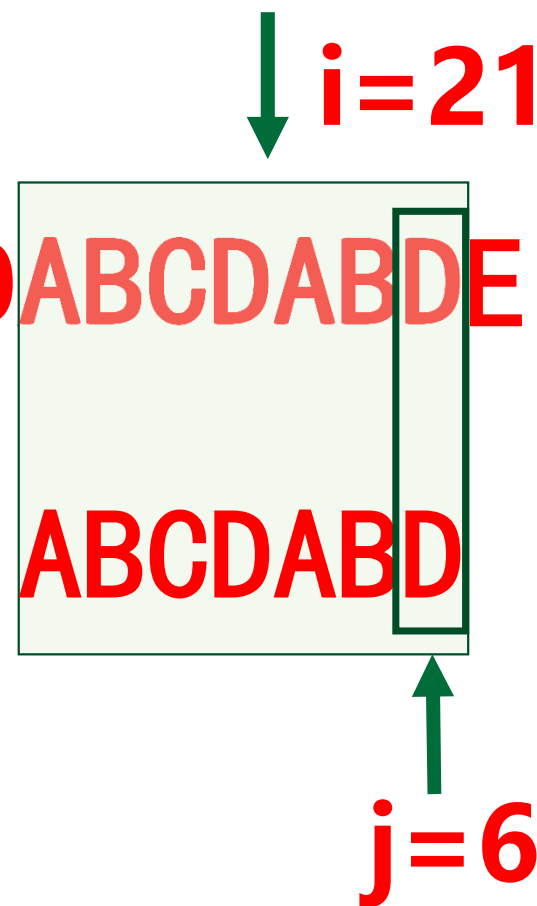
模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

串的模式匹配算法：KMP算法

串S: BBC ABCDAB ABCDABCDABDE

串P:

i始终没有回溯，一直向右，提高了效率



模式串	A	B	C	D	A	B	D
最大长度值	0	0	0	0	1	2	0
next 数组	-1	0	0	0	0	1	2

KMP改进方法

KMP改进方法

串S: **abacababc**



串T: **abab**

	a	b	a	b
索引值	0	1	2	3
最大前缀后缀	0	0	1	2
next[j]	-1	0	0	1
是否满足	初始值 无须优化	$P[1] \neq P[\text{next}[1]]$	$P[2] == P[\text{next}[2]]$	$p[3] == p[\text{next}[3]]$
是否优化	初始值 无须优化	不优化	需要优化	需要优化
优化的next[j]	-1	0	$\text{next}[2] = \text{next}[\text{next}[2]] = \text{next}[0] = -1$	$\text{next}[3] = \text{next}[\text{next}[3]] = \text{next}[1] = 0$

KMP改进方法

串S: **abacababc**



串T: **abab**

	a	b	a	b
索引值	0	1	2	3
最大前缀后缀	0	0	1	2
next[j]	-1	0	0	1
优化的next[j]	-1	0	-1	0

KMP改进方法

串S: **abacababc**



串T: **abab**

重复比较

串S: **abacababc**



abab

按照next值，应该移到1号位置

	a	b	a	b
索引值	0	1	2	3
最大前缀后缀	0	0	1	2
next[j]	-1	0	0	1
优化的next[j]	-1	0	-1	0

KMP改进方法

串S: **abacababc**

串T: **abab**

串S: **abacababc**

abab

按照优化的next值，应该移到0号位置

	a	b	a	b
索引值	0	1	2	3
最大前缀后缀	0	0	1	2
next[j]	-1	0	0	1
优化的next[j]	-1	0	-1	0

KMP改进方法

```
int j,k;
j=0;k=-1;
nextval[0]=-1;
while(j<S.length-1){
    if (k==-1 || S[j]==S[k] ){
        j++;k++;
        if (S[j])!=S[k]{ nextval[j]=k; }
        else nextval[j]=nextval[k];
    }
    else
        k=nextval[k];
}
```

Boyer-Moore算法 (BM算法)

KMP的匹配是从模式串的开头开始匹配的，而1977年，德克萨斯大学的Robert S. Boyer教授和J Strother Moore教授发明了一种新的字符串匹配算法：**Boyer-Moore算法，简称BM算法**。该算法从模式串的尾部开始匹配，且拥有在最坏情况下 $O(N)$ 的时间复杂度。在实践中，比KMP算法的实际效能高。

Boyer-Moore算法 (BM算法)

BM算法定义了两个规则：

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。

好后缀规则：当字符失配时，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，且如果好后缀在模式串中没有再次出现，则为-1。

Boyer-Moore算法 (BM算法)

串S: HERE IS A SIMPLE EXAMPLE

串P: EXAMPLE

"文本串"与"模式串"头部对齐，从尾部开始比较。"S"与"E"不匹配。这时，"S"就被称为"坏字符" (bad character)，即不匹配的字符，它对应着模式串的第6位。且"S"不包含在模式串"EXAMPLE"之中（相当于最右出现位置是-1），这意味着可以把模式串后移 $6 - (-1) = 7$ 位，从而直接移到"S"的后一位。

Boyer-Moore算法 (BM算法)

串S: HERE IS A SIMPLE EXAMPLE

串P: EXAMPLE

从尾部开始比较。“P”与“E”不匹配，坏字符，它对应着模式串的第6位。但“P”包含在模式串“EXAMPLE”之中（出现位置是4），这意味着可以把模式串后移 $6-4=2$ 位，从而直接后移2位，两个“P”对齐。

Boyer-Moore算法 (BM算法)

串S: HERE IS A SIMPLE EXAMPLE

串P: EXAMPLE

“P” 与 “P” 匹配。这时，“P” 就被称为 “好字符”，即所有尾部匹配的字符串。注意，“MPLE”、“PLE”、“LE”、“E”都是好后缀。

Boyer-Moore算法 (BM算法)

串S: HERE IS A SIMPLE EXAMPLE

串P:

EXAMPLE

“I” 与 “A” 不匹配。这时，按照坏字符规则，向右移3位($2 - (-1)$)。
所有的 “好字符”，即所有尾部匹配的字符串。注意，“MPLE”、
“PLE”、“LE”、“E” 都是好后缀。后移位数 = 好后缀在模式串中的
位置 - 好后缀在模式串中上一次出现的位置，且如果好后缀在模式串中没
有再次出现，则为-1。所有的 “好后缀” (MPLE、PLE、LE、E) 之中，
只有 “E” 在 “EXAMPLE” 的尾部和头部出现，所以后移 $6 - 0 = 6$ 位。

Boyer-Moore算法 (BM算法)

串S: HERE IS A SIMPLE EXAMPLE

串P:

EXAMPLE

Boyer-Moore算法 (BM算法)

串S: HERE IS A SIMPLE EXAMPLE

串P: EXAMPLE

“P”与“E”不匹配。这时，“P”就被称为“坏字符” (bad character)，它对应着模式串的第6位。且“P”包含在模式串“EXAMPLE”之中（最右出现位置是4），这意味着可以把模式串后移 $6 - 4 = 2$ 位。

Boyer-Moore算法 (BM算法)

串S: HERE IS A SIMPLE EXAMPLE

串P: EXAMPLE

“E” 与 “E” 匹配。这时，“E” 就被称为 “好字符” (good character)，它对应着模式串的第6位。“EXAMPLE” “XAMPLE” “AMPLE” “MPLE”、“PLE”、“LE”、“E” 都是好后缀。

Sunday算法

BM算法虽然通常比KMP算法快，但BM算法也还不是现有字符串查找算法中最快的算法，最后再介绍一种比BM算法更快的查找算法即**Sunday算法**。

Sunday算法由Daniel M.Sunday在1990年提出，它的思想跟BM算法很相似：

只不过Sunday算法是从前往后匹配，在匹配失败时关注的是文本串中参加匹配的最末位字符的下一位字符。

如果该字符没有在模式串中出现则直接跳过，即移动位数 = 匹配串长度 + 1；

否则，其移动位数 = 模式串中最右端的该字符到末尾的距离 + 1。

串S: substrⁱing searching algorithm

串P: search

结果发现在第2个字符处发现不匹配，不匹配时关注文本串中参加匹配的最末位字符的下一位字符，即标粗的字符 i ，因为模式串 `search` 中并不存在 i ，所以模式串直接跳过一大片，向右移动位数 = 匹配串长度 + 1 = $6 + 1 = 7$ ，从 i 之后的那个字符（即字符 n ）开始下一步的匹配。

串S: substring searching algorithm

串P: search

第一个字符就不匹配，再看文本串中参加匹配的最末位字符的下一位字符，是'r'，它出现在模式串中的倒数第3位，于是把模式串向右移动3位（r 到模式串末尾的距离 + 1 = 2 + 1 = 3），使两个'r'对齐

串S: substring searching algorithm

串P:

search

匹配成功

本章学习要点

1. 掌握串类型的特点，并能在相应的应用问题中正确选用它们。
2. 熟练KMP实现方法
3. 了解BM算法和Sunday算法