串 (Sequence)

@M了个J

https://github.com/CoderMJLee http://cnblogs.com/mjios

KMP详解: https://www.zhihu.com/question/21923021



码拉松





\$\$ \$\bar{\text{NON-BOTTOMS}}\$ \$\bar{\text{Sequence}}\$\$

■ 本课程研究的串是开发中非常熟悉的字符串,是由若干个字符组成的有限序列

String text = "Thank";						
Т	h	a	n	k		

■字符串 thank 的前缀 (prefix)、真前缀 (proper prefix)、后缀 (suffix)、真后缀 (proper suffix)

前缀	t, th, tha, than, thank
真前缀	t, th, tha, than
后缀	thank, hank, ank, nk, k
真后缀	hank, ank, nk, k



MAR 中四四算法

- 本课程主要研究串的匹配问题,比如
- □查找一个模式串 (pattern) 在文本串 (text) 中的位置

```
String text = "Hello World";
String pattern = "or";
text.indexOf(pattern); // 7
text.indexOf("other"); // -1
```

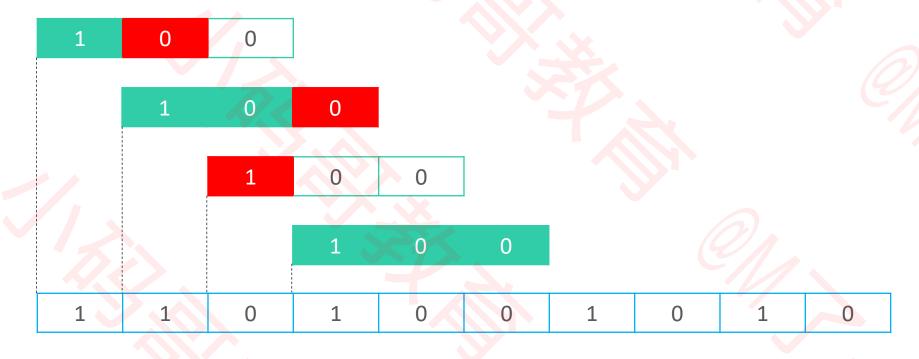
- ■几个经典的串匹配算法
- □蛮力 (Brute Force)
- **□** KMP
- Boyer-Moore
- Karp-Rabin
- **□** Sunday

■本课程用 tlen 代表文本串 text 的长度, plen 代表模式串 pattern 的长度



本の日本教育 蛮力 (Brute Force)

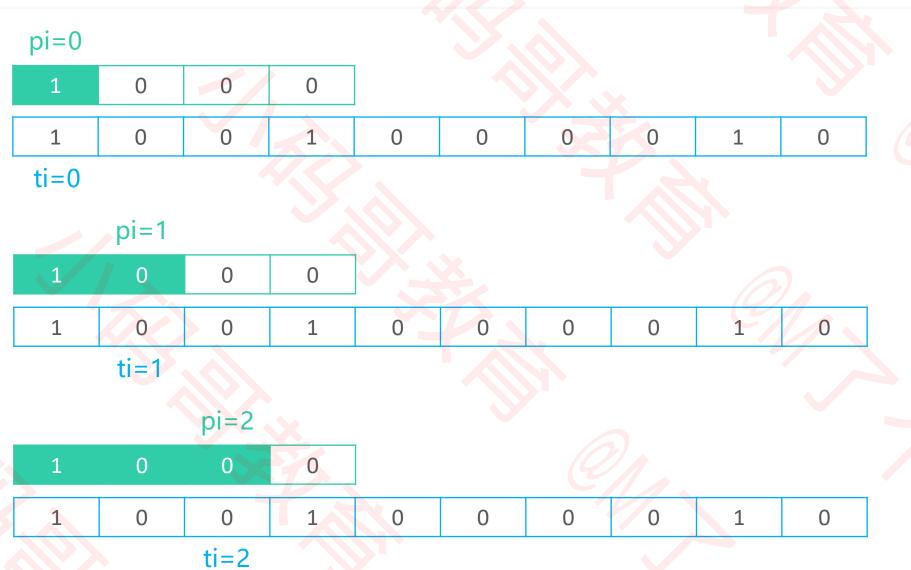
■ 以字符为单位,从左到右移动模式串,直到匹配成功



■ 蛮力算法有 2 种常见实现思路



小码母教育 SEEMYGO **蛋力1 — 执行过程**

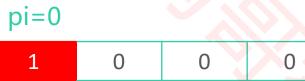


■ pi 的取值范围 [0, plen)

■ ti 的取值范围 [0, tlen)

小四哥教育 盛力1 — 执行过程





1	0	0	1	0	0 0	0	1 0

		3	1	0	0	0		
1	0	0	1	0	0	0 0	1	0
						ti=7		

$$pi = 0$$

 $ti = pi - 1$

pi == plen 代表匹配成功

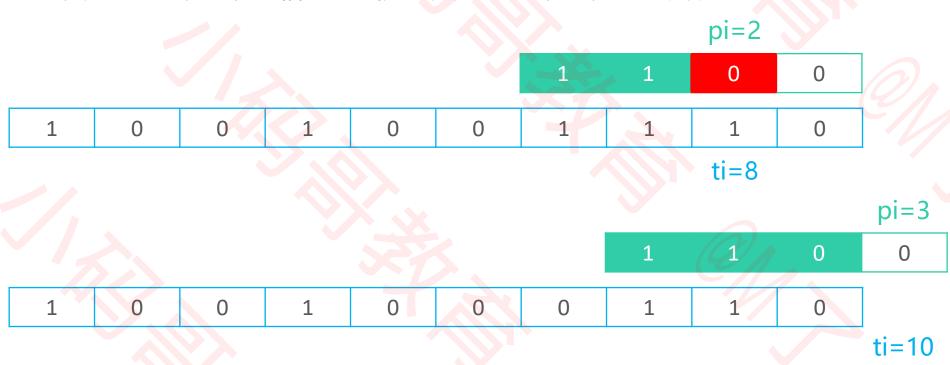
小码哥教育 SEEMYGO 重力1一实现

```
public static int indexOf(String text, String pattern) {
   if (text == null || pattern == null) return -1;
   int tlen = text.length();
   int plen = pattern.length();
   if (tlen == 0 || plen == 0 || tlen < plen) return -1;
   int pi = 0, ti = 0;
   while (pi < plen && ti < tlen) {</pre>
       if (text.charAt(ti) == pattern.charAt(pi)) {
            ti++;
            pi++;
        } else {
            ti -= pi - 1;
            pi = 0;
   return pi == plen ? ti - pi : -1;
```



小码哥教育 SEEMYGO **蛮力1**一优化

■ 此前实现的蛮力算法,在恰当的时候可以提前退出,减少比较次数



■ 因此, ti 的退出条件可以从 ti < tlen 改为

□ti – pi 是指每一轮比较中 text 首个比较字符的位置

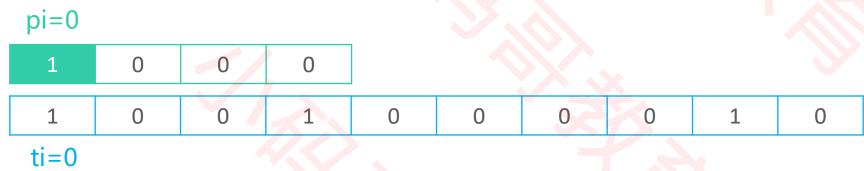
这是完全没必要的比较

小門哥教育 蛮力1一优化实现

```
public static int indexOf(String text, String pattern) {
    if (text == null || pattern == null) return -1;
    int tlen = text.length();
    int plen = pattern.length();
   if (tlen == 0 || plen == 0 || tlen < plen) return -1;
   int pi = 0, ti = 0;
   int tmax = tlen - plen;
   while (pi < plen && ti - pi <= tmax) {</pre>
        if (text.charAt(ti) == pattern.charAt(pi)) {
            ti++;
            pi++;
        } else {
            ti -= pi - 1;
            pi = 0;
    return pi == plen ? ti - pi : -1;
```



↑ 小妈哥教育 **蛮力2 – 执行过程**



pi=1

1	0	0	0

1 0 0 1	0 0	0 1 0

ti + pi ti=0

1 0	0	0						
1 0	0	1	0	0	0	0	1	0

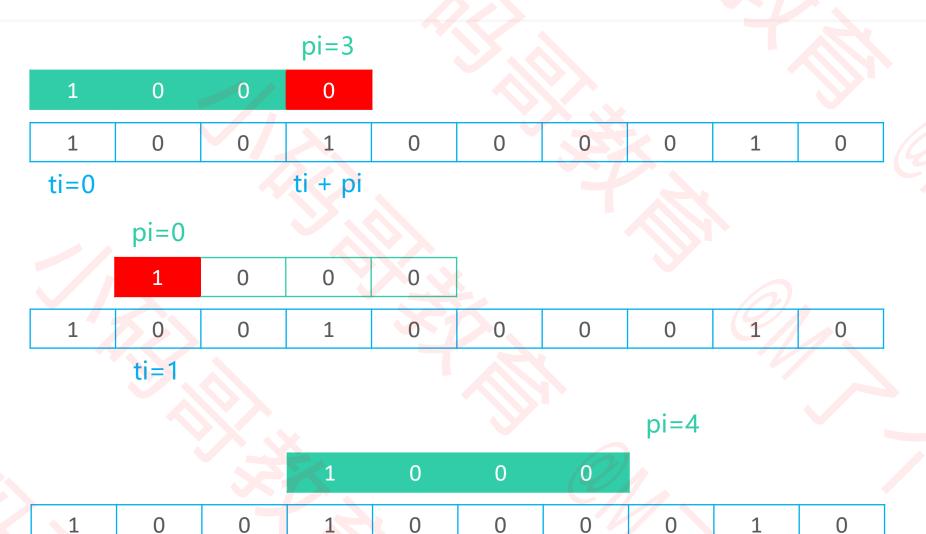
■ pi 的取值范围 [0, plen)

■ ti 的取值范围 [0, tlen – plen]



小門司教育 SEEMYGO 蛮力2 – 执行过程

ti=3



ti + pi

pi = 0ti++

> pi == plen 代表匹配成功

小码哥教育 SEEMYGO 事力2一实现

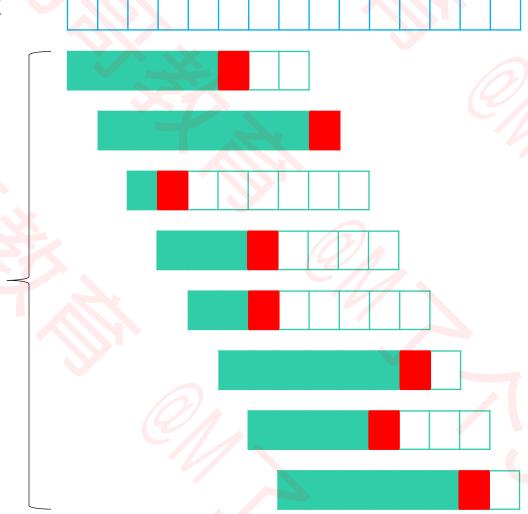
```
public static int indexOf(String text, String pattern) {
    if (text == null || pattern == null) return -1;
   int tlen = text.length();
    int plen = pattern.length();
    if (tlen == 0 || plen == 0 || tlen < plen) return -1;
    int tmax = tlen - plen;
    for (int ti = 0; ti <= tmax; ti++) {</pre>
        int pi = 0;
        for (; pi < plen; pi++) {
            if (text.charAt(ti + pi) != pattern.charAt(pi)) break;
        if (pi == plen) return ti;
    return -1;
```



小小門哥教育 蛮力 — 性能分析

■ n 是文本串长度, m 是模式串长度

最多 n - m + 1 轮





MAR A MAR

- ■最好情况
- □只需一轮比较就完全匹配成功, 比较 m 次 (m 是模式串的长度)
- □时间复杂度为 O(m)

1	0			
1	0	0	1	0

- 最坏情况 (字符集越大, 出现概率越低)
- □执行了 n-m+1 轮比较 (n是文本串的长度)
- □每轮都比较至模式串的末字符后失败 (m-1次成功,1次失败)
- □时间复杂度为 O(m * (n m + 1)), 由于一般 m 远小于 n, 所以为 O(mn)

						1	1	0	
1	1	1	1	1	1	1	1	1	1



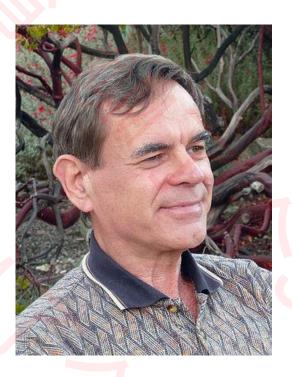
■ KMP 是 Knuth-Morris-Pratt 的简称(取名自3位发明人的名字),于1977年发布



Donald Knuth



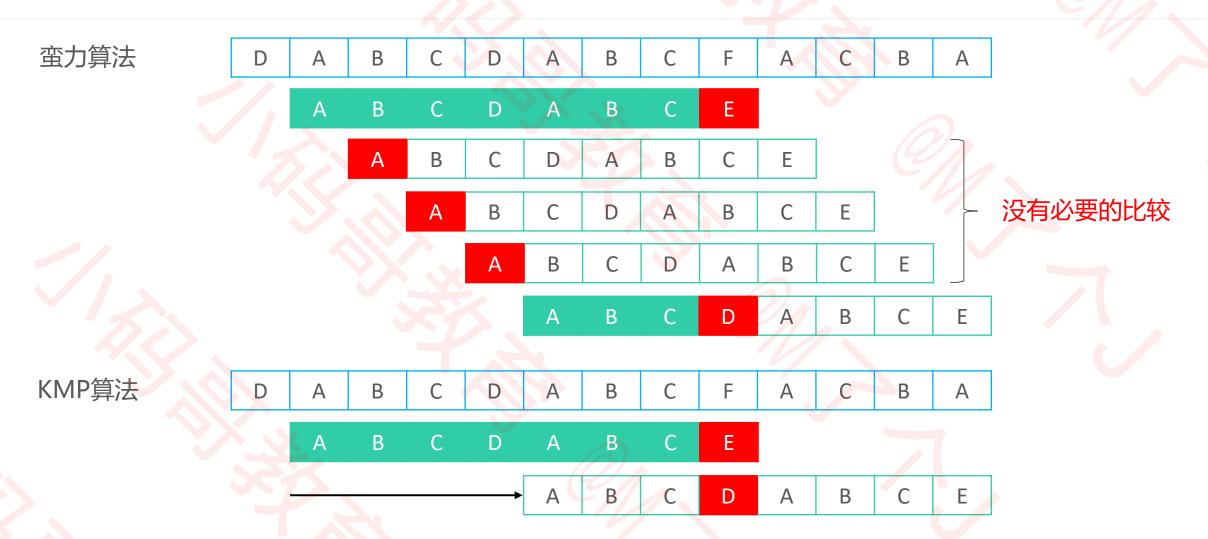
James Hiram Morris



Vaughan Pratt



小码哥教育 蛋力 VS KMP



■ 对比蛮力算法,KMP的精妙之处: 充分利用了此前比较过的内容, 可以很聪明地跳过一些不必要的比较位置

小門司教育 KMP — next表的使用

■ KMP 会预先根据模式串的内容生成一张 next 表 (一般是个数组)

	模	式串"	ABCDAI	BCE"的	next	表		
模式串字符	А	В	С	D	Α	В	С	Е
索引	0	1	2	3	4	5	6	7
元素	-1	0	0	0	0	1	2	3

ti=8



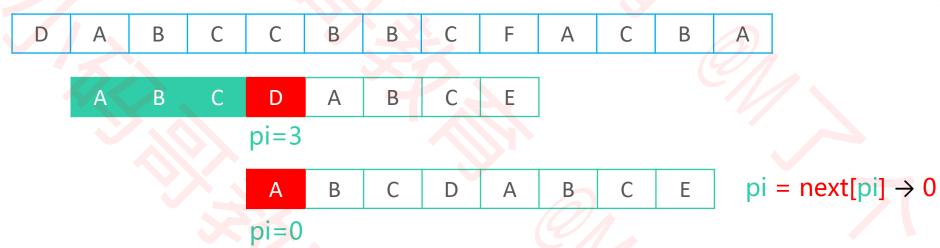
$$pi = next[pi] \rightarrow 3$$



小码 哥教育 KMP — next表的使用

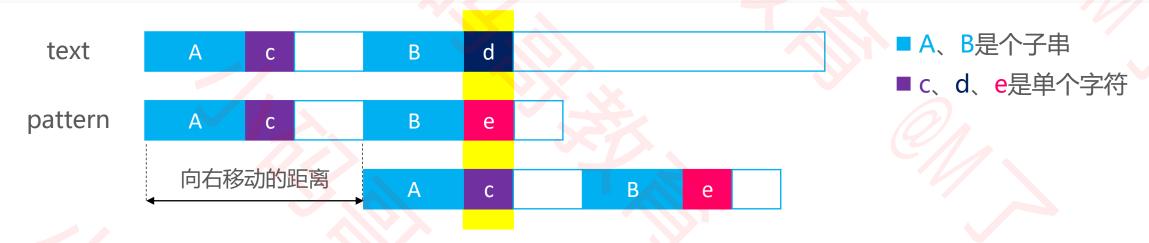
	模式串 "ABCDABCE" 的 next 表							
模式串字符	A	В	С	D	А	В	С	Е
索引	0	1	2	3	4	5	6	7
元素	-1	0	0	0	0	1	2	3

$$ti=5$$





小码 引教育 KMP 一核心原理



- 当 d、e 失配时,如果希望 pattern 能够一次性向右移动一大段距离,然后直接比较 d、c 字符
- □前提条件是 A 必须等于 B
- 所以 KMP 必须在失配字符 e 左边的子串中找出符合条件的 A、B,从而得知向右移动的距离
- 向右移动的距离: e左边子串的长度 A的长度, 等价于: e的索引 c的索引
- 且 c的索引 == next[e的索引], 所以向右移动的距离: e的索引 next[e的索引]
- ■总结
- □如果在 pi 位置失配,向右移动的距离是 pi next[pi],所以 next[pi] 越小,移动距离越大
- □ next[pi] 是 pi 左边子串的真前缀后缀的最大公共子串长度



增量 KMP - 真前缀后缀的最大公共子串长度

模式串	真前缀	真后缀	最大公共子串长度
ABCDABCE	A, AB, ABC, ABCD, ABCDA, ABCDAB, ABCDABC	BCDABCE, CDABCE, DABCE, ABCE, BCE, CE, E	0
ABCDABC	A, AB, ABC, ABCD, ABCDA, ABCDAB	BCDABC, CDABC, DABC, ABC, BC, C	3
ABCDAB	A, AB, ABC, ABCD, ABCDA	BCDAB, CDAB, DAB, AB, B	2
ABCDA	A, AB, ABC, ABCD	BCDA, CDA, DA, A	1
ABCD	A, AB, ABC	BCD, CD, D	0
ABC	A, AB	BC, C	0
AB	A > 7	В	0
А			0

模式串字符	Α	В	C	D	Α	В	С	Е
最大公共子串长度	0	0	0	0	1	2	3	0



Number of the N

模式串字符	А	В	С	D	А	В	C	E
最大公共子串长度	0	0	0	0	1	2	3	0

■ 将最大公共子串长度都向后移动 1 位,首字符设置为 负1,就得到了 next 表

模式串 "ABCDABCE" 的 next 表									
模式串字符	Α	В	С	D	Α	В	С	E	
索引	0	1	2	3	4	5	6	70	
元素	-1	0	0	0	0	1	2	3	



имя КМР — 负1的精妙之处



小四哥教育 KMP - 主算法实现

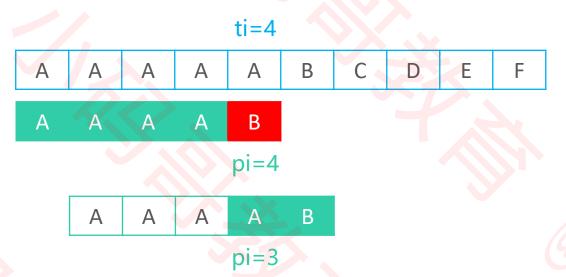
```
public static int indexOf(String text, String pattern) {
    if (text == null || pattern == null) return -1;
    int plen = pattern.length();
    int tlen = text.length();
    if (tlen == 0 \mid \mid plen == 0 \mid \mid tlen < plen) return <math>-1;
    int[] next = next(pattern);
    int pi = 0, ti = 0;
    int tmax = tlen - plen;
    while (pi < plen && ti - pi <= tmax) {</pre>
        if (pi < 0 || text.charAt(ti) == pattern.charAt(pi)) {</pre>
             ti++;
             pi++;
        } else {
             pi = next[pi];
    return pi == plen ? ti - pi : -1;
```



增量 KMP - 为什么是"最大"公共子串长度?

■ 假设文本串是AAAAABCDEF,模式串是AAAAB

模式串	真前缀	真后缀	公共子串长度
AAAA	A, AA, AAA	A, AA, AAA	1, 2, 3
AAA	A, AA	A, AA	1, 2
AA	A	A	1



pi=1

- 应该将1、2、3中的哪个值赋值给 pi 是正确的?
- 将 3 赋值给 pi
- □向右移动了1个字符单位,最后成功匹配
- 将 1 赋值给 pi
- □向右移动了3个字符单位,错过了成功匹配的机会
- 公共子串长度越小,向右移动的距离越大,越不安全
- 公共子串长度越大,向右移动的距离越小,越安全

小門司教育 KMP — next表的构造思路



- 已知 next[i] == n
- ① 如果 pattern.charAt(i) == pattern.charAt(n)
- □那么 next[i + 1] == n + 1
- ② 如果 pattern.charAt(i)!= pattern.charAt(n)
- □已知 next[n] == k
- ■如果 pattern.charAt(i) == pattern.charAt(k)
- ✓ 那么 next[i + 1] == k + 1
- □如果 pattern.charAt(i)!= pattern.charAt(k)
- ✓ 将 k 代入 n , 重复执行 ②

小門司教育 KMP — next表的代码实现

```
public static int[] next(String pattern) {
    int len = pattern.length();
    int[] next = new int[len];
    int i = 0;
    int n = next[i] = -1;
    int imax = len - 1;
    while (i < imax) {</pre>
        if (n < 0 | pattern.charAt(i) == pattern.charAt(n)) {</pre>
            next[++i] = ++n;
        } else {
            n = next[n];
    return next;
```



小門司教息 KMP - next表的不足之处

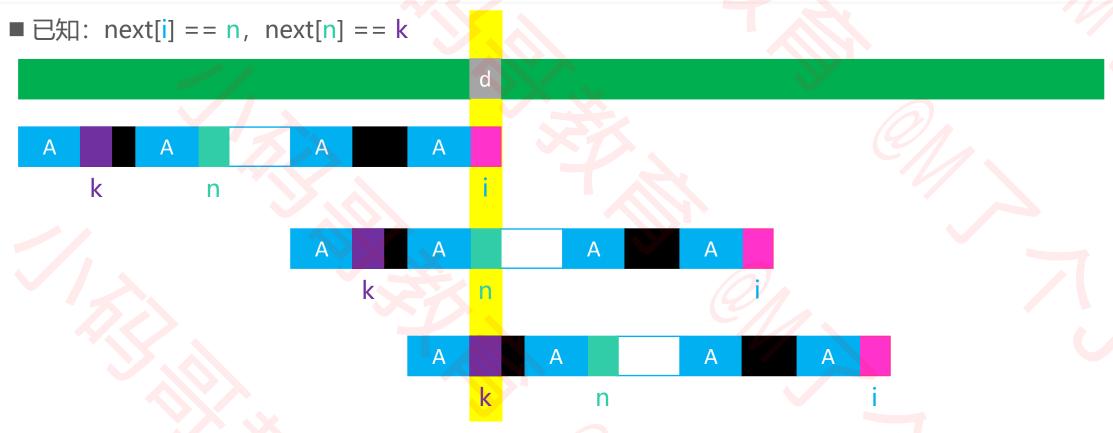
■ 假设文本串是 AAABAAAAB ,模式串是 AAAAB



■ 在这种情况下,KMP显得比较笨拙



小門司教息 KMP — next表的优化思路



- 如果 pattern[i]!= d, 就让模式串滑动到 next[i] (也就是n) 位置跟 d 进行比较
- 如果 pattern[n]!= d, 就让模式串滑动到 next[n] (也就是k) 位置跟 d 进行比较
- ■如果 pattern[i] == pattern[n],那么当 i 位置失配时,模式串最终必然会滑到 k 位置跟 d 进行比较
- ■所以 next[i] 直接存储 next[n] (也就是k) 即可



小門司教息 KMP - next表的优化实现

```
public static int[] next(String pattern) {
    int len = pattern.length();
    int[] next = new int[len];
    int i = 0;
    int n = next[i] = -1;
    int imax = len - 1;
    while (i < imax) {</pre>
        if (n < 0 | pattern.charAt(i) == pattern.charAt(n)) {</pre>
            i++;
            n++;
            if (pattern.charAt(i) == pattern.charAt(n)) {
                next[i] = next[n];
            } else {
                next[i] = n;
        } else {
            n = next[n];
    return next;
```



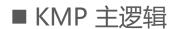
Myseemys。 KMP - next表的优化效果

模式串 "AAAAB" 的 next 表									
模式串字符	Α	Α	Α	Α	В				
索引	0	1	2	3	4				
优化前	-1	0	1	2	3				
优化后	-1	-1	-1	-1	3				

A	Α	А	В	А	А	А	A	В
А	A	A	Α	В				
				А	Α	Α	Α	В



小码 哥教育 KMP — 性能分析



□最好时间复杂度: 0(m)

□最坏时间复杂度: 0(n), 不超过0(2n)

■ next 表的构造过程跟 KMP 主体逻辑类似

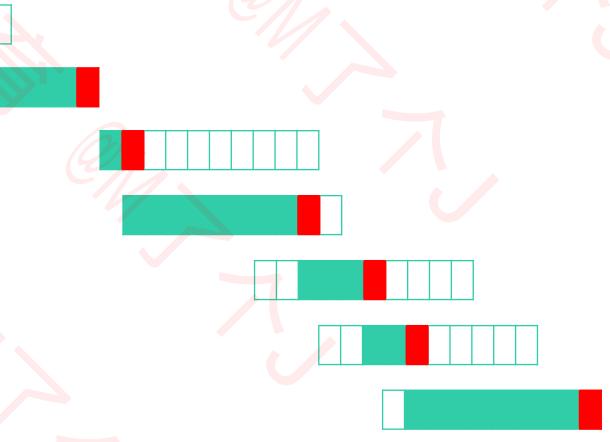
□时间复杂度: 0(m)

■ KMP 整体

□最好时间复杂度: O(m)

□最坏时间复杂度: O(n+m)

□空间复杂度: O(m)





Myga 重力 vs KMP

- 蛮力算法为何低效?
- ■当字符失配时
- □蛮力算法
- ✓ ti 回溯到左边位置
- ✓ pi 回溯到 0
- ■KMP 算法
- ✓ ti 不必回溯
- ✓ pi 不一定要回溯到 0