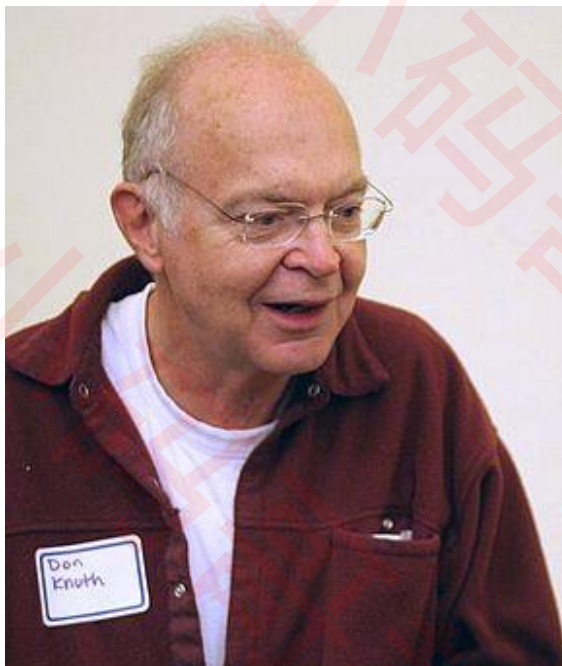


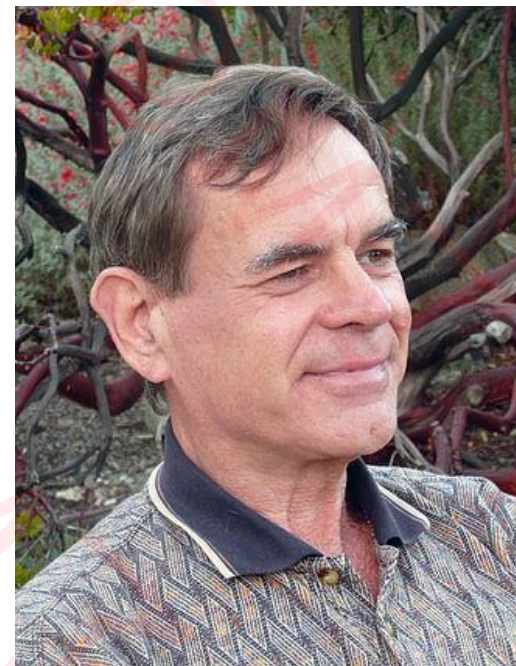
■ KMP 是 **K**nuth-**M**orris-**P**ratt 的简称（取名自3位发明人的名字），于1977年发布



Donald **K**nuth



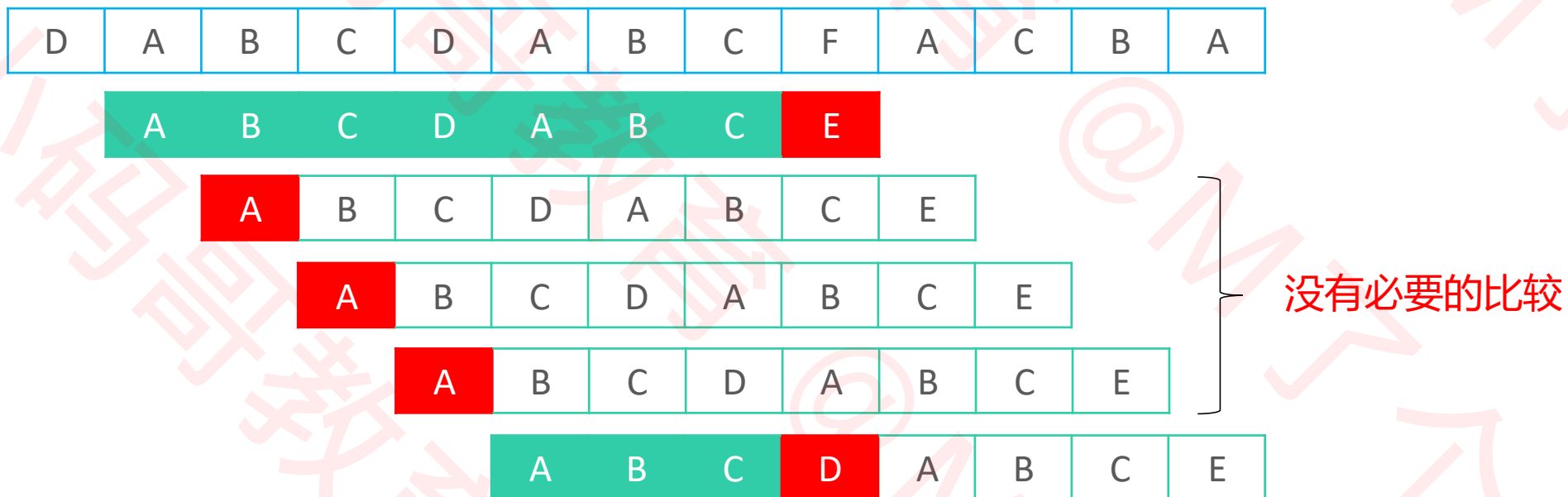
James Hiram **M**orris



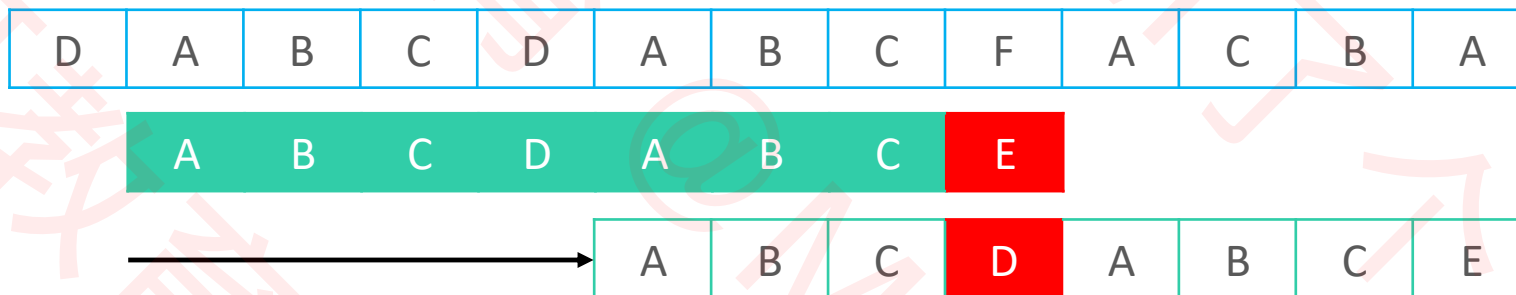
Vaughan **P**ratt



# 蛮力算法



## KMP算法



■ 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

# 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 "ABCDABCE" 的 next 表								
模式串字符	A	B	C	D	A	B	C	E
索引	0	1	2	3	4	5	6	7
元素	-1	0	0	0	0	1	2	3

$ti=8$

D	A	B	C	D	A	B	C	F	A	C	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---

A	B	C	D	A	B	C	E
---	---	---	---	---	---	---	---

$pi - next[pi] \rightarrow 4$

$pi=7$

向右移动的距离

A	B	C	D	A	B	C	E
---	---	---	---	---	---	---	---

$pi=3$

$pi = next[pi] \rightarrow 3$

# KMP – next表的使用

模式串 "ABCDABCE" 的 next 表								
模式串字符	A	B	C	D	A	B	C	E
索引	0	1	2	3	4	5	6	7
元素	-1	0	0	0	0	1	2	3

ti=5

D	A	B	C	C	B	B	C	F	A	C	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---

A	B	C	D	A	B	C	E
---	---	---	---	---	---	---	---

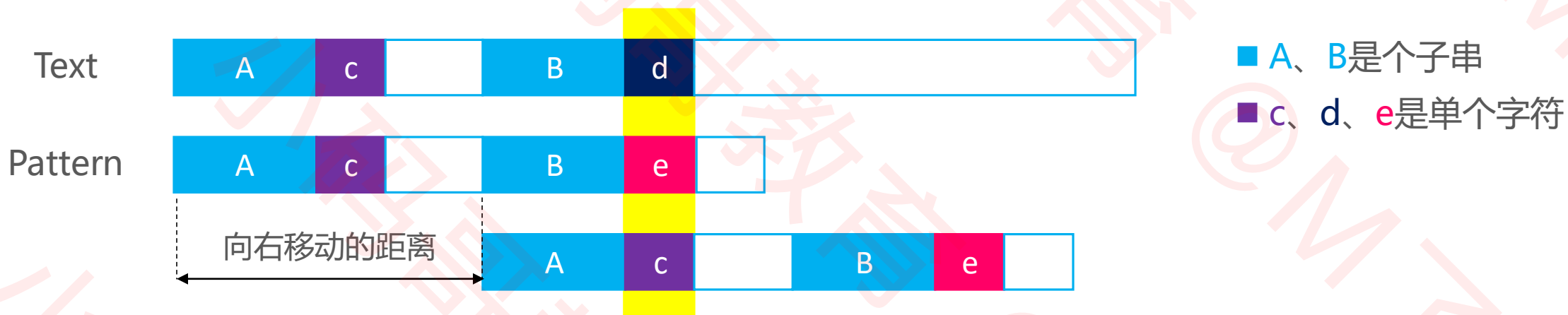
pi=3

A	B	C	D	A	B	C	E
---	---	---	---	---	---	---	---

pi=0

pi = next[pi] → 0

# KMP – 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

# 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, <b>ABC</b> , ABCD, ABCDA, ABCDAB	BCDABC, CDABC, DABC, <b>ABC</b> , BC, C	3
ABCDAB	A, <b>AB</b> , ABC, ABCD, ABCDA	BCDAB, CDAB, DAB, <b>AB</b> , B	2
ABCD A	<b>A</b> , AB, ABC, ABCD	BCDA, CDA, DA, <b>A</b>	1
ABCD	A, AB, ABC	BCD, CD, D	0
ABC	A, AB	BC, C	0
AB	A	B	0
A			0

模式串字符	A	B	C	D	A	B	C	E
最大公共子串长度	0	0	0	0	1	2	3	0

# KMP – 得到next表

模式串字符	A	B	C	D	A	B	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	A	B	C	D	A	B	C	E
索引	0	1	2	3	4	5	6	7
元素	-1	0	0	0	0	1	2	3

# KMP – 负1的精妙之处

ti=2

D	A	C	C	D	A	B	C	F	A	C	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---

A	B	C	D	A	B	C	E
---	---	---	---	---	---	---	---

pi=0

ti=3

D	A	C	C	D	A	B	C	F	A	C	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---

A	B	C	D	A	B	C	E
---	---	---	---	---	---	---	---

pi=0

ti=2

D	A	C	C	D	A	B	C	F	A	C	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---

*	A	B	C	D	A	B	C	E
---	---	---	---	---	---	---	---	---

pi=-1

pi = next[pi] → -1

pi++ → 0

ti++ → 3

■ 相当于在负1位置有个假想的通配字符（哨兵）

□ 匹配成功后 ti++、pi++



# 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 || plen == 0 || tlen < plen) return -1;  
    int[] next = next(pattern);  
    int pi = 0, ti = 0;  
    int tmax = tlen - plen;  
    while (pi < plen && ti - pi <= tmax) {  
        if (pi < 0 || text.charAt(ti) == pattern.charAt(pi)) {  
            ti++;  
            pi++;  
        } else {  
            pi = next[pi];  
        }  
    }  
    return pi == plen ? ti - pi : -1;  
}
```

# 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

ti=4

A	A	A	A	A	B	C	D	E	F
---	---	---	---	---	---	---	---	---	---

A	A	A	A	B
---	---	---	---	---

pi=4

A	A	A	A	B
---	---	---	---	---

pi=3

A	A	A	A	B
---	---	---	---	---

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

# KMP – next表的构造思路



■ 已知  $\text{next}[i] == n$

① 如果  $\text{Pattern}[i] == \text{Pattern}[n]$

□ 那么  $\text{next}[i + 1] == n + 1$

② 如果  $\text{Pattern}[i] \neq \text{Pattern}[n]$

□ 已知  $\text{next}[n] == k$

□ 如果  $\text{Pattern}[i] == \text{Pattern}[k]$

✓ 那么  $\text{next}[i + 1] == k + 1$

□ 如果  $\text{Pattern}[i] \neq \text{Pattern}[k]$

✓ 将  $k$  代入  $n$ ，重复执行 ②

# 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) {  
        if (n < 0 || pattern.charAt(i) == pattern.charAt(n)) {  
            next[++i] = ++n;  
        } else {  
            n = next[n];  
        }  
    }  
    return next;  
}
```

# KMP – next表的不足之处

■ 假设文本串是 AAABAAAAB，模式串是 AAAAB

模式串 "AAAAB" 的 next 表					
模式串字符	A		A	A	B
索引	0	1	2	3	4
元素	-1	0	1	2	3

A A A B A A A A B

A A A A B

A A A A B

A A A A B

A A A A B

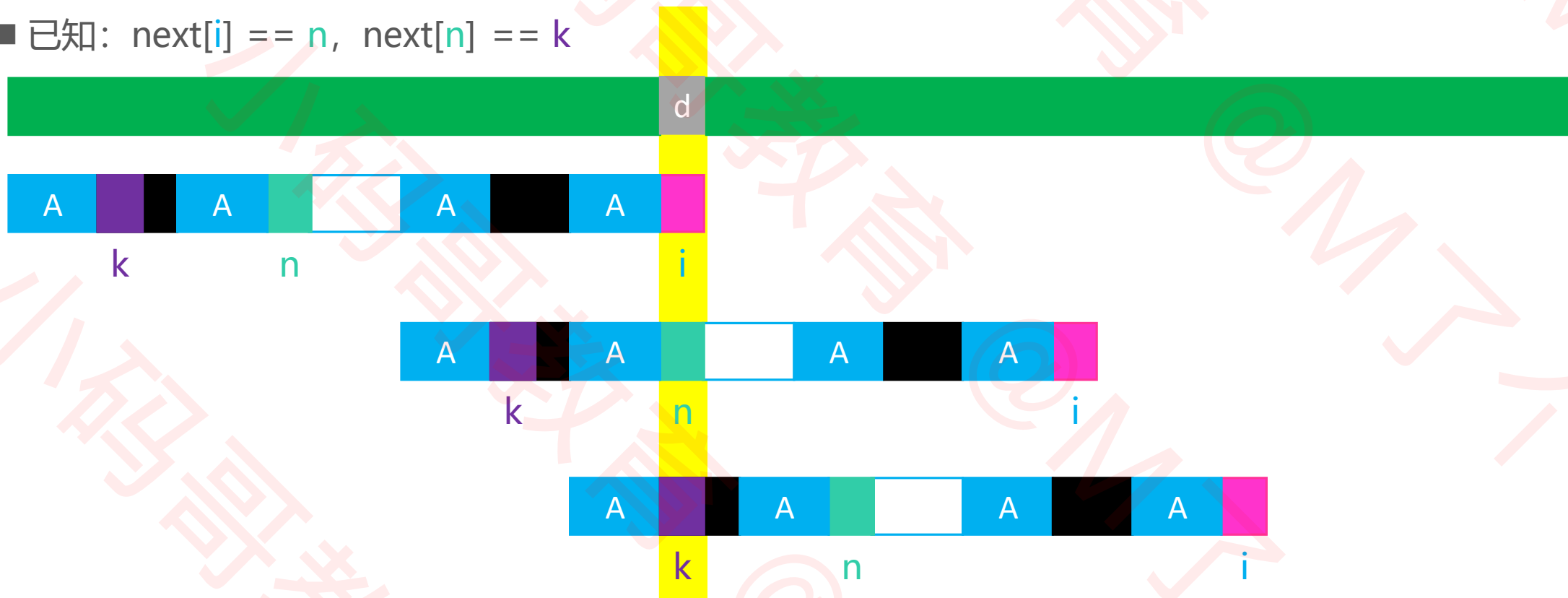
A A A A B

没有必要的比较

■ 在这种情况下，KMP显得比较笨拙

# KMP – next表的优化思路

■ 已知:  $\text{next}[i] == n$ ,  $\text{next}[n] == k$



- 如果  $\text{Pattern}[i] \neq d$ , 就让模式串滑动到  $\text{next}[i]$  (也就是n) 位置跟 d 进行比较
- 如果  $\text{Pattern}[n] \neq d$ , 就让模式串滑动到  $\text{next}[n]$  (也就是k) 位置跟 d 进行比较
- 如果  $\text{Pattern}[i] == \text{Pattern}[n]$ , 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  $\text{next}[i]$  直接存储  $\text{next}[n]$  (也就是k) 即可

# 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) {  
        if (n < 0 || pattern.charAt(i) == pattern.charAt(n)) {  
            i++;  
            n++;  
            if (pattern.charAt(i) == pattern.charAt(n)) {  
                next[i] = next[n];  
            } else {  
                next[i] = n;  
            }  
        } else {  
            n = next[n];  
        }  
    }  
    return next;  
}
```

# KMP – next表的优化效果

模式串 "AAAAB" 的 next 表

模式串字符	A	A	A	A	B
索引	0	1	2	3	4
优化前	-1	0	1	2	3
优化后	-1	-1	-1	-1	3

A	A	A	B	A	A	A	A	B
---	---	---	---	---	---	---	---	---

A	A	A	A	B
---	---	---	---	---

A	A	A	A	B
---	---	---	---	---



# KMP – 性能分析

## ■ KMP 主逻辑

□ 最好时间复杂度:  $O(m)$

□ 最坏时间复杂度:  $O(n)$ , 不超过  $O(2n)$

## ■ next 表的构造过程跟 KMP 主体逻辑类似

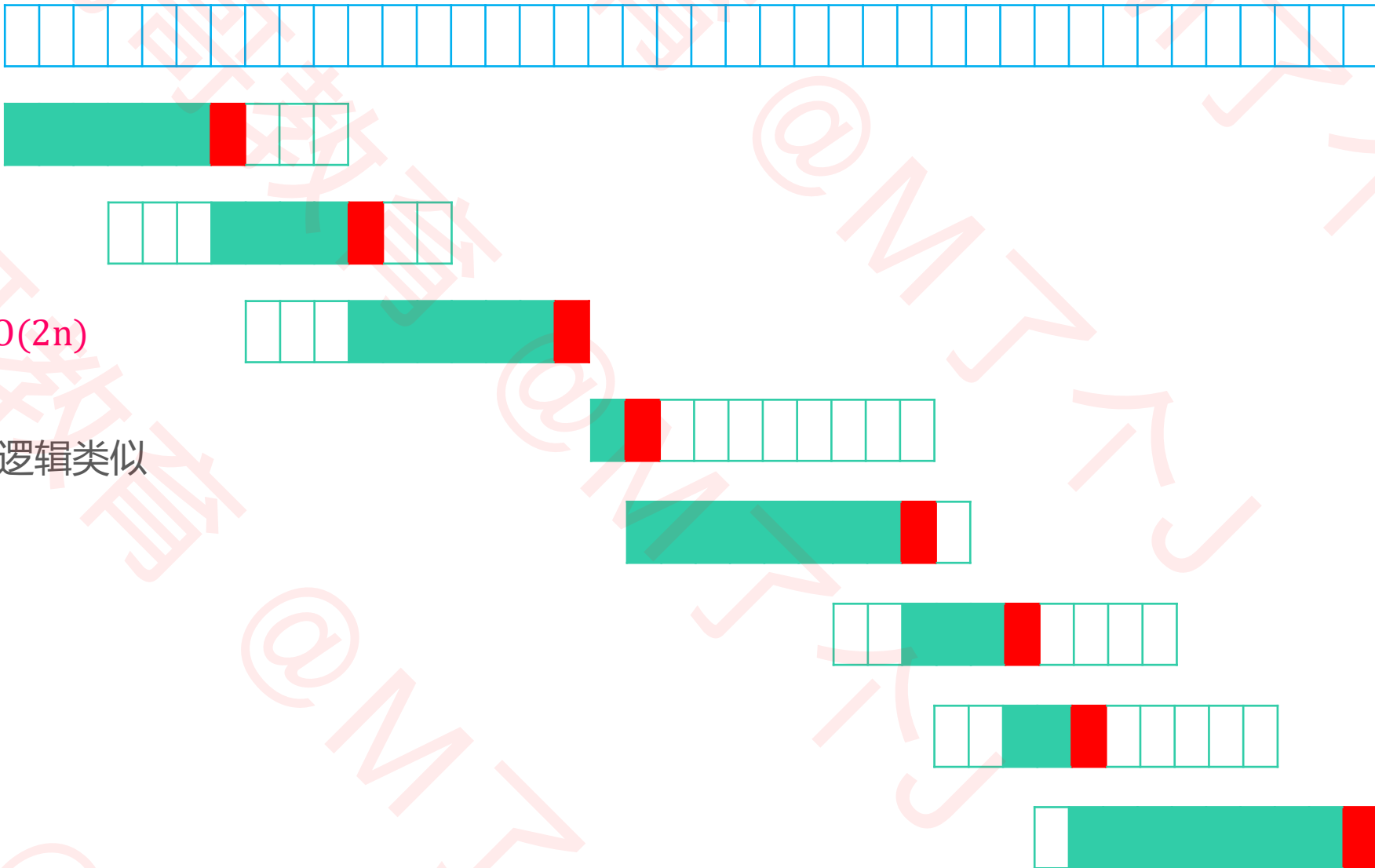
□ 时间复杂度:  $O(m)$

## ■ 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:  $O(m)$



# 蛮力 vs KMP

- 蛮力算法为何低效?

- 当字符失配时

- 蛮力算法

- ✓  $t_i$  回溯到左边位置

- ✓  $p_i$  回溯到 0

- KMP 算法

- ✓  $t_i$  不必回溯

- ✓  $p_i$  不一定要回溯到 0