
Jheng-Huang Hong

Department of Computer Science and Information Engineering
National Cheng Kung University
Tainan, Taiwan



String Basic

- 字串 string
 - 字元的有序序列 $A = a_0a_1 \dots a_{n-1}$
 - 字元集元集是字串的長度
- 子字串 substring
 - $(A[i, j]) = a_i a_{i+1} a_{i+2} \dots a_j$ (A 連續的一段)
- 子序列 subsequence
 - $B = a_{q_1} a_{q_2} a_{q_3} \dots a_{q_m}, 0 \leq q_1 < q_2 < \dots < q_m < n$ (不連續)
- 後綴 suffix
 - A 的一個子字串 $S_A(k) = a_k a_{k+1} a_{k+2} \dots a_n, 0 \leq k < n$
- 前綴 prefix
 - A 的一個子字串 $P_A(h) = a_0 a_1 a_2 \dots a_h, 0 \leq h < n$



String Basic

- S = “abcbbab”
 - 子字串： “ bcb” , “ bba” , ...
 - 子序列： “ acb” , “ bbb” , ...
 - 前綴： “ abcb” , “ ab” , ...
 - 後綴： “ bbab” , “ ab” , ...



String Matching

~	!	@	#	\$	%	^	&	*	()	-	=	Delete
Tab	Q	W	E	R	T	Y	U	I	O	P	{	}	
Caps	A	S	D	F	G	H	J	K	L	:	"	'	Enter
Shift	Z	X	C	V	B	N	M	<	>	?	/	Shift	
Ctrl		Alt									Alt		Ctrl

- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

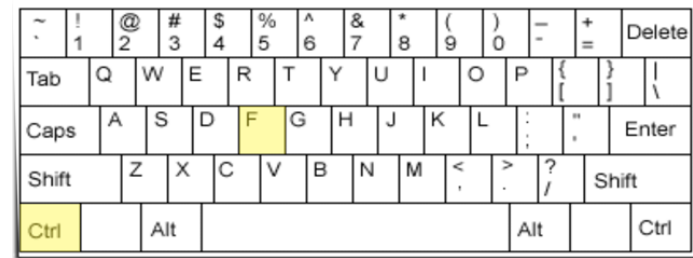
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "abcde~~fg~~"
- B = "cde"



String Matching

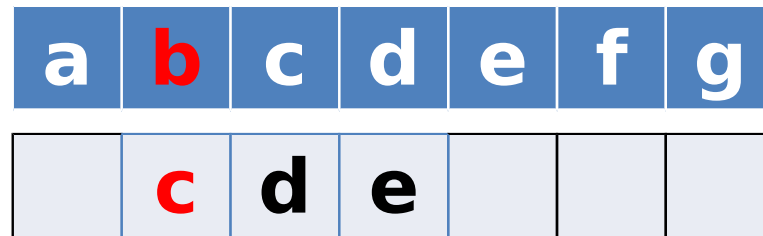


- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

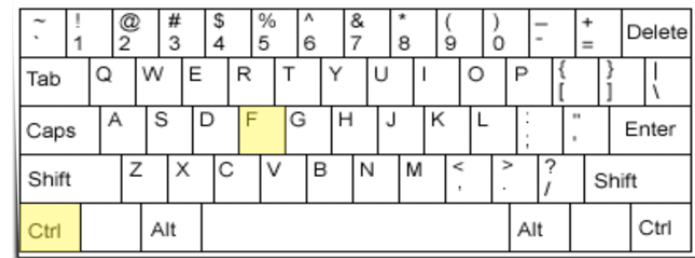
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "abcde~~fg~~"
- B = "cde"



String Matching



- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

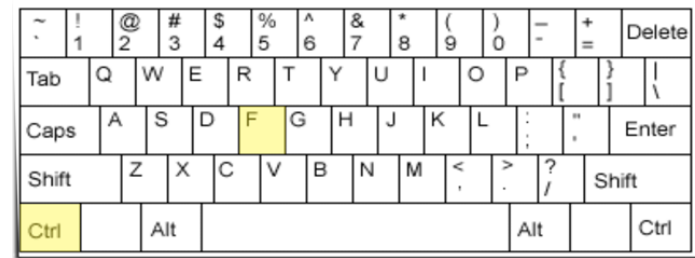
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "abcde~~fg~~"
- B = "cde"



String Matching

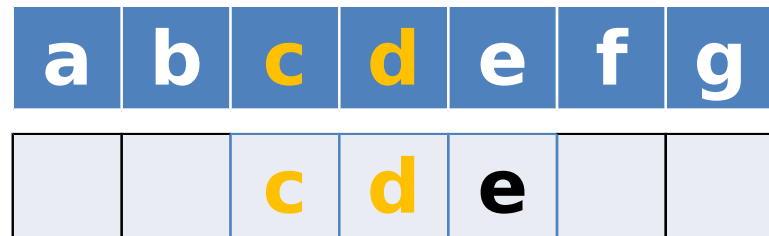


- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

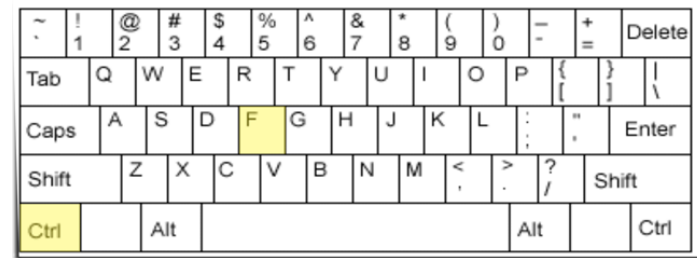
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "abcde~~fg~~"
- B = "cde"



String Matching



- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

Matchin
g!!

- 複雜度: $O(|A|)$
- A = "abcde~~fg~~"
- B = "cde"



String Matching

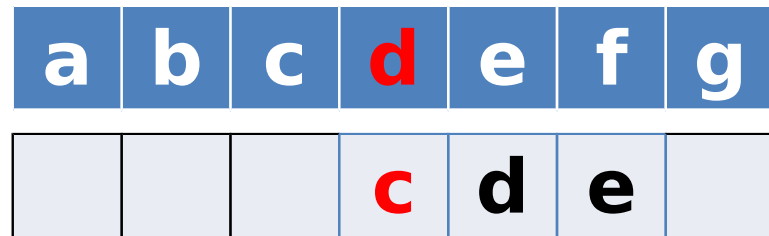
~	!	@	#	\$	%	^	&	*	()	-	=	Delete
Tab	Q	W	E	R	T	Y	U	I	O	P	{	}	
Caps	A	S	D	F	G	H	J	K	L	:	"	'	Enter
Shift	Z	X	C	V	B	N	M	<	>	?	/	Shift	
Ctrl		Alt									Alt		Ctrl

- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "abcde~~fg~~"
- B = "cde"



String Matching

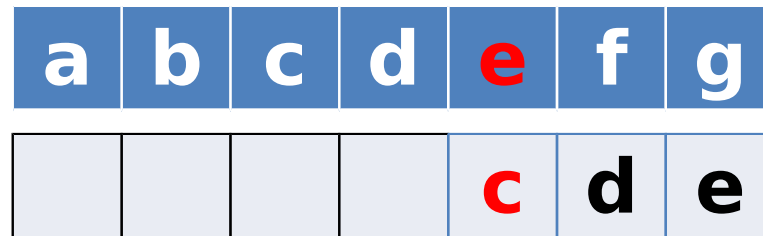
~	!	@	#	\$	%	^	&	*	()	-	=	Delete
Tab	Q	W	E	R	T	Y	U	I	O	P	{	}	
Caps	A	S	D	F	G	H	J	K	L	:	"	'	Enter
Shift	Z	X	C	V	B	N	M	<	>	?	/	Shift	
Ctrl		Alt									Alt		Ctrl

- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "abcde~~fg~~"
- B = "cde"



String Matching

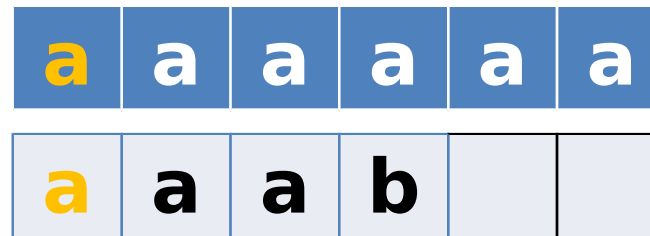
~	!	@	#	\$	%	^	&	*	()	-	=	Delete
Tab	Q	W	E	R	T	Y	U	I	O	P	{	}	
Caps	A	S	D	F	G	H	J	K	L	:	"	'	Enter
Shift	Z	X	C	V	B	N	M	<	>	?	/	Shift	
Ctrl		Alt									Alt		Ctrl

- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

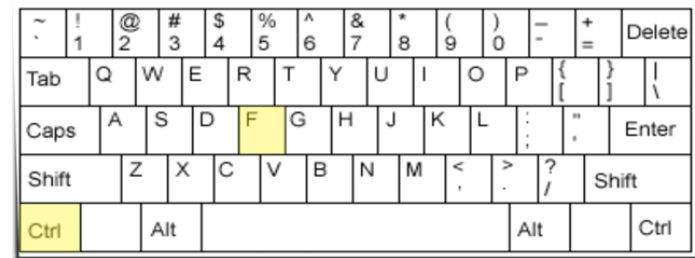
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "aaaaaaaa::aaa"
- B = "aaaaaaaa::aab"



String Matching

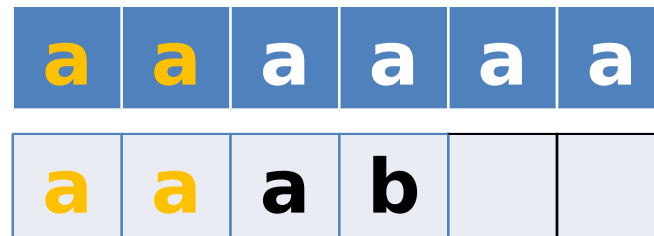


- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

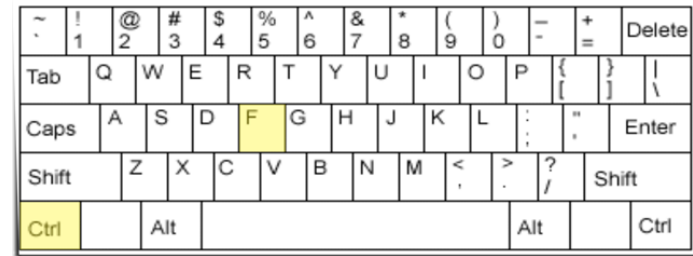
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "aaaaaaa::aaa"
- B = "aaaaaaa::aab"



String Matching

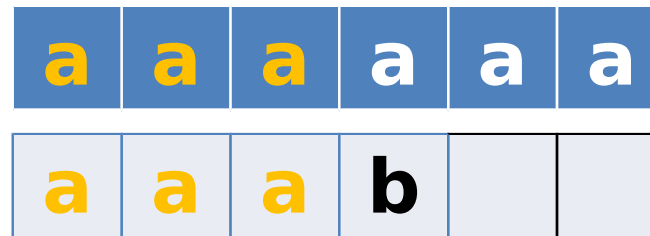


- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "aaaaaaaa::aaa"
- B = "aaaaaaaa::aab"



String Matching

~	!	@	#	\$	%	^	&	*	()	-	=	Delete
Tab	Q	W	E	R	T	Y	U	I	O	P	{	}	
Caps	A	S	D	F	G	H	J	K	L	:	"	'	Enter
Shift	Z	X	C	V	B	N	M	<	>	?	/	Shift	
Ctrl		Alt									Alt		Ctrl

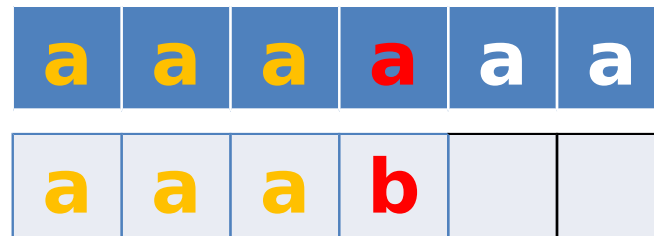
- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

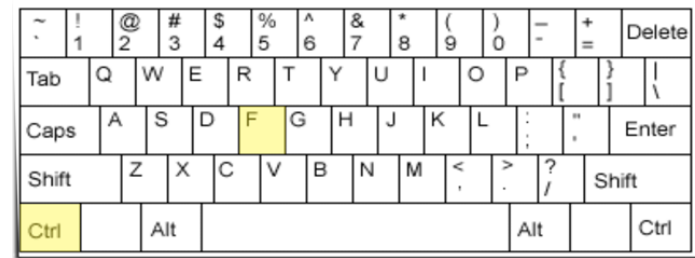
```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$

- A = "aaaaaaaa::aaa"
- B = "aaaaaaaa::aab"



String Matching



- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

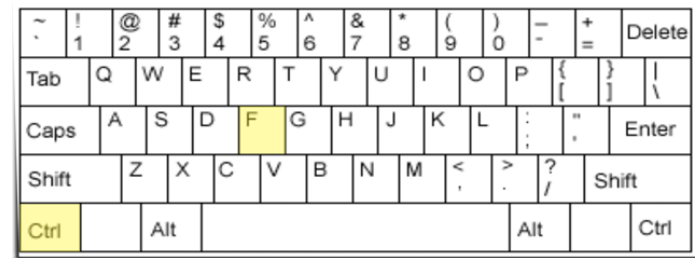
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "aaaaaaaa::aaa"
- B = "aaaaaaaa::aab"



String Matching

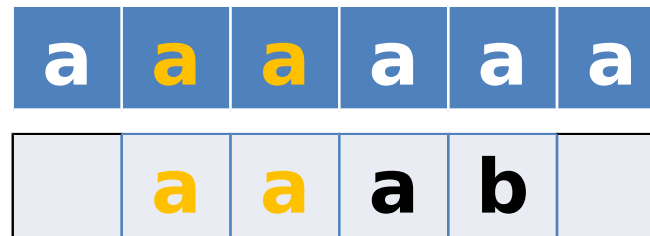


- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

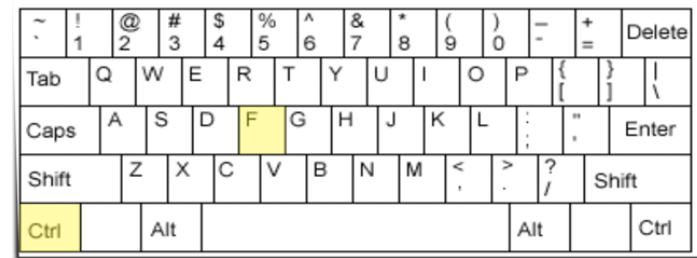
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- $A = \text{"aaaaaa::aaa"}$
- $B = \text{"aaaaaa::aab"}$



String Matching



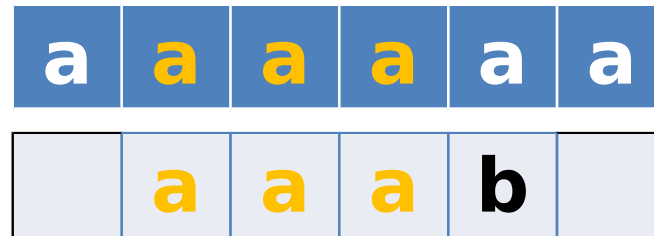
- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

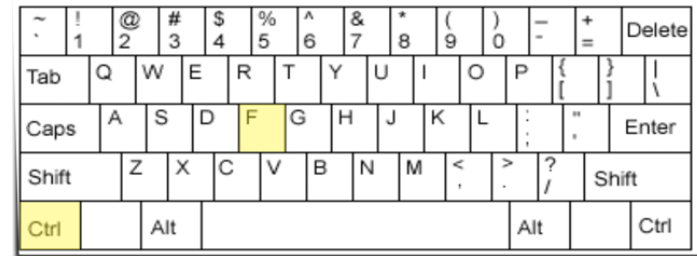
```

1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
    
```

- 複雜度： $O(|A|)$
- $A = \text{"aaaaaaaa::aaa"}$
- $B = \text{"aaaaaaaa::aab"}$



String Matching

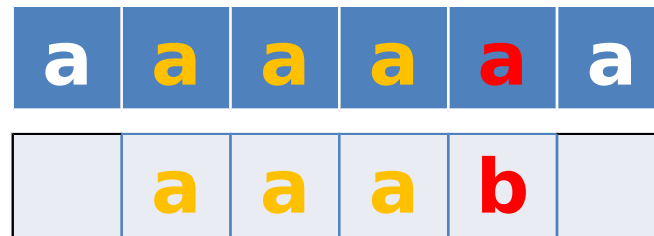


- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- $A = \text{"aaaaaaaa::aaa"}$
- $B = \text{"aaaaaaaa::aab"}$



String Matching

~	!	@	#	\$	%	^	&	*	()	-	=	Delete
Tab	Q	W	E	R	T	Y	U	I	O	P	{	}	
Caps	A	S	D	F	G	H	J	K	L	:	"	'	Enter
Shift	Z	X	C	V	B	N	M	<	>	?	/	Shift	
Ctrl		Alt									Alt		Ctrl

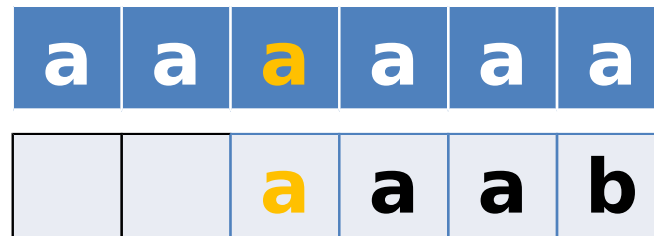
- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$

- A = "aaaaaaaa::aaa"
- B = "aaaaaaaa::aab"



String Matching

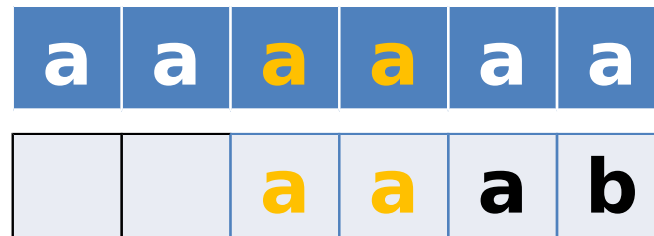
~	!	@	#	\$	%	^	&	*	()	-	=	Delete
Tab	Q	W	E	R	T	Y	U	I	O	P	{	}	
Caps	A	S	D	F	G	H	J	K	L	:	"	'	Enter
Shift	Z	X	C	V	B	N	M	<	>	?	/	Shift	
Ctrl		Alt									Alt		Ctrl

- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

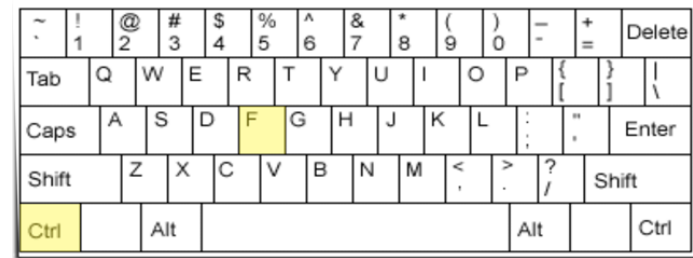
•

```
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
```

- 複雜度： $O(|A|)$
- A = "aaaaaaaa::aaa"
- B = "aaaaaaaa::aab"



String Matching



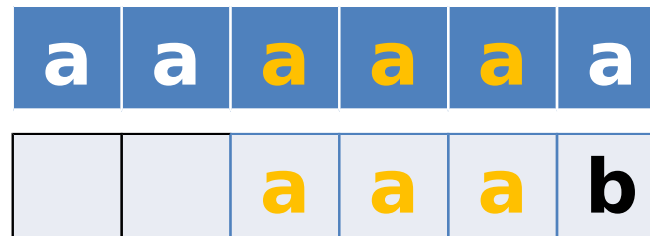
- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

•

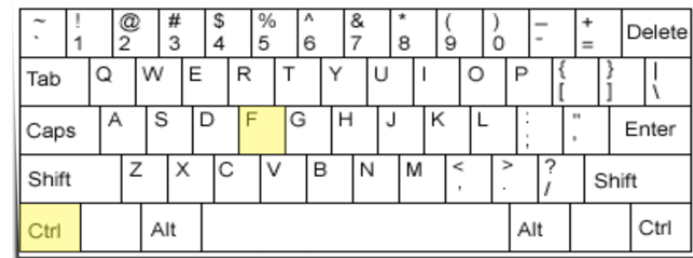
```

1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
    
```

- 複雜度： $O(|A|)$
- $A = \text{"aaaaaa::aaa"}$
- $B = \text{"aaaaaa::aab"}$



String Matching



- 給兩個字串 A, B 找出所有 B 出現在 A 中的位置

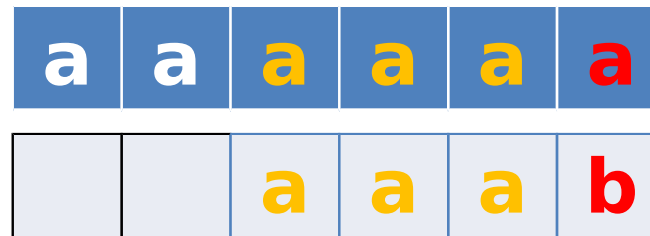
•

```

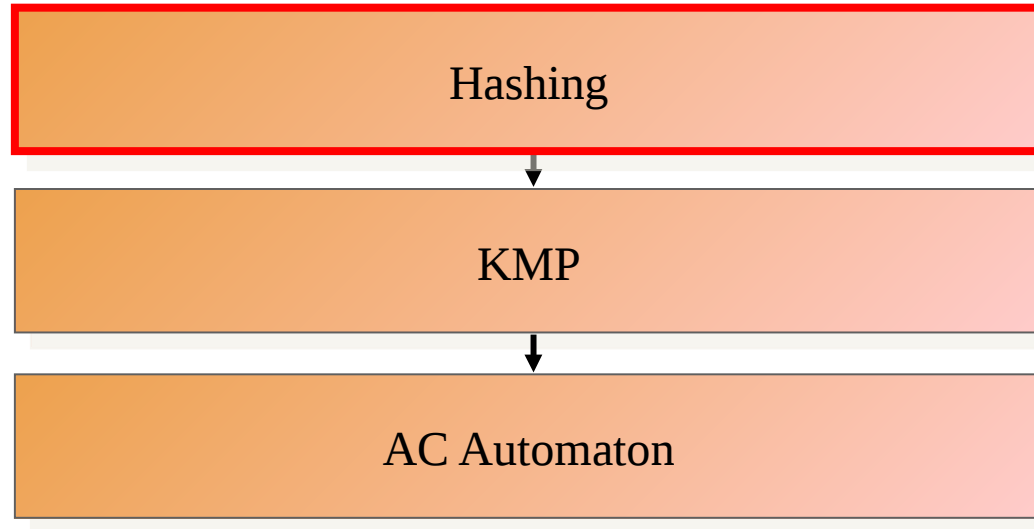
1 for(int i=0; i+lenB<=lenA; ++i){
2     int mat=0;
3     while(mat<lenB && A[i+mat]==B[mat]) ++mat;
4     if(mat == lenB) print(i);
5 }
    
```

- 複雜度： ~~$O(|A|)$~~ $O(|A||B|)$

- A = "aaaaaaaa::aaa"
- B = "aaaaaaaa::aab"



Outline



Hashing

- 分類

 - 將字串分到有限的整數裡

 - 函數 $f: string \mapsto \{0, 1, \dots, Q - 1\}$

- 要求

 - 容易取得

 - 均勻

- 思考

1. $f(A) \neq f(B) \Rightarrow A \neq B$

2. $A \neq B \Rightarrow f(A) \neq f(B)$ → 不一定

3. 分 n 類；碰撞機率 $1/n$



Hashing

- Rabin-Karp rolling hash function 定義

- $f(A) = a_0p^{n-1} + a_1p^{n-2} + \dots + a_{n-2}p + a_{n-1} \text{ mod } q$

- 類似: p 進位制, 分成 q 類

- p, q 取不同質數 \Rightarrow 均勻

- 滾動

1. $f(A) \equiv f(A[0, n-2])p + a_{n-1} \text{ mod } q$

- \Rightarrow 計算 A 所有前綴的 hash value $O(|A|)$

2. $f(A[i, j]) \equiv f(A[0, j]) - p^{j-i+1} f(A[0, i-1]) \text{ mod } q$

- \Rightarrow 任何 A 子字串的 hash value $O(1)$

3. 枚舉 A 長度為 l 的子字串, 比較 hash value

- $\Rightarrow O(N)$



Hashing

- $A = \text{"abcdefg"}$

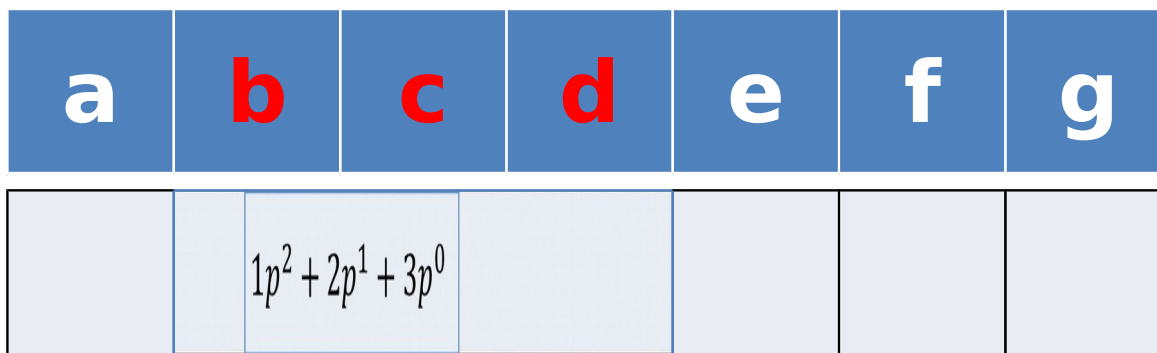
a	b	c	d	e	f	g
$0p^2 + 1p^1 + 2p^0$						

- $B = \text{"cde"}$ $= 2p^2 + 3p^1 + 4p^0$



Hashing

- $A = \text{"abcdefg"}$

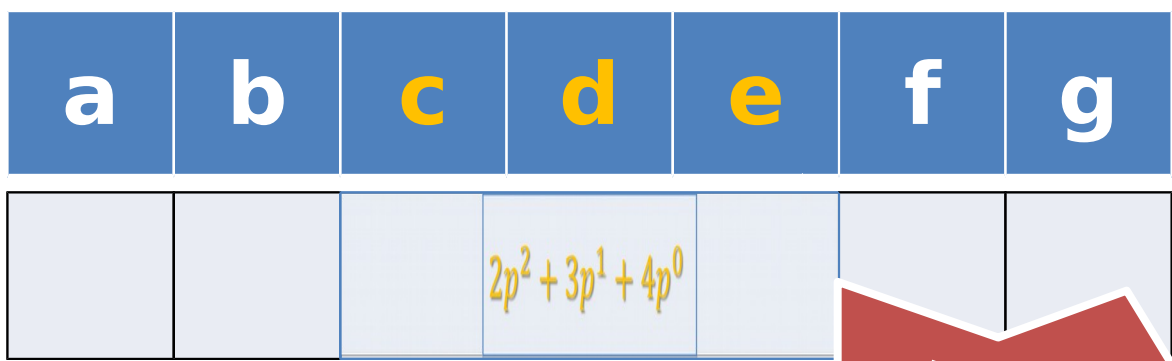


- $B = \text{"cde"}$ $= 2p^2 + 3p^1 + 4p^0$



Hashing

- $A = \text{"abcdefg"}$

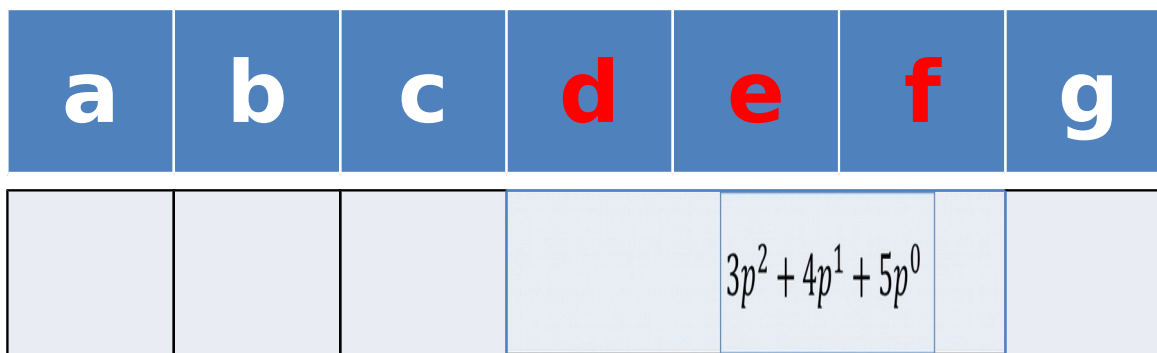


- $B = \text{"cde"}$ = $2p^2 + 3p^1 + 4p^0$



Hashing

- $A = \text{"abcdefg"}$

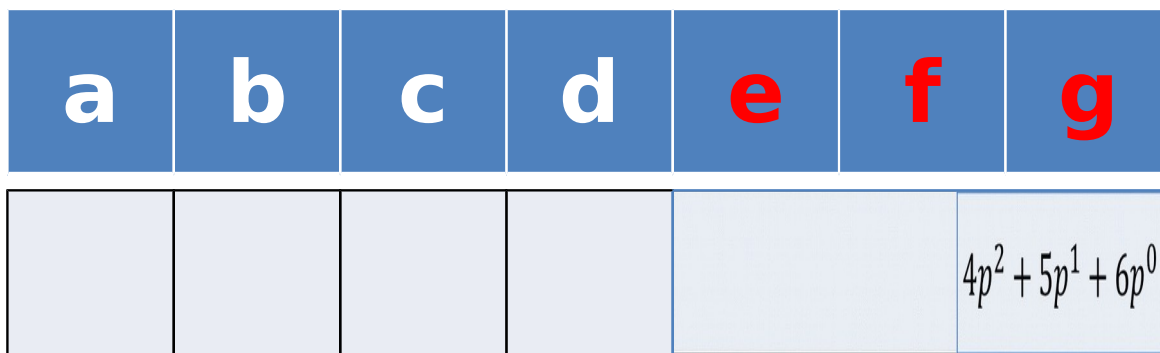


- $B = \text{"cde"}$ $= 2p^2 + 3p^1 + 4p^0$



Hashing

- $A = \text{"abcdefg"}$



- $B = \text{"cde"}$ $= 2p^2 + 3p^1 + 4p^0$



Hashing

- 回測剛剛

- $A \neq B \Rightarrow f(A) \neq f(B)$ → 不一定

- 相等時重新檢查一次?

- $A = "aaaaaaaa::aaa"$
 - $B = "aaaaaaaa::aab"$

- q 取大一點 (long long 質數)

- 碰撞機率小

- ex. $q \in 10^{15} \Rightarrow probability: 10^{-15}$

- ex. 2147483647



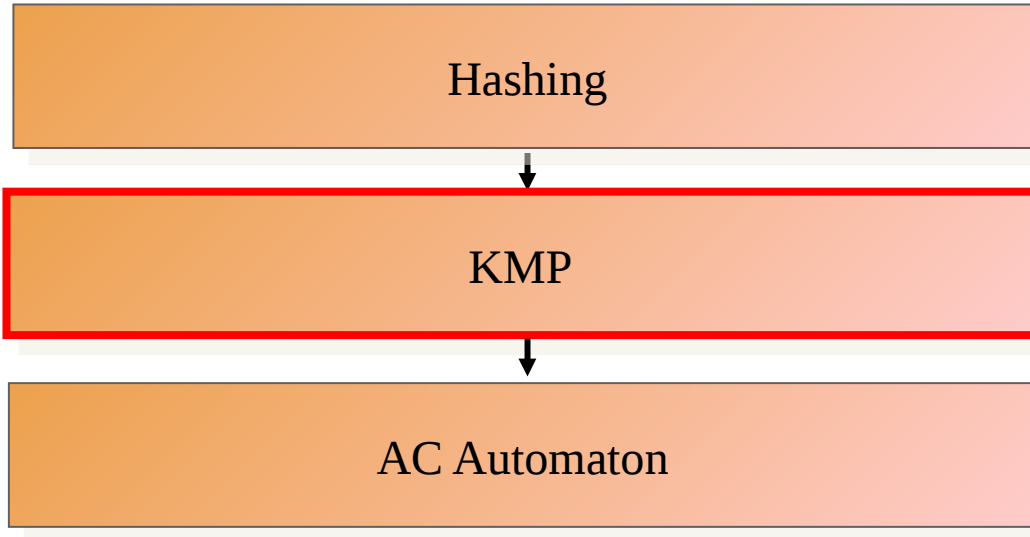
Hashing – 參考

```
1 #define MAXN 1000000
2 #define prime_mod 1073676287
3 /*prime_mod 必須要是質數*/
4 typedef long long T;
5 char s[MAXN+5];
6 T h[MAXN+5];/*hash陣列*/
7 T h_base[MAXN+5];/*h_base[n]=(prime^n)%prime_mod*/
8 inline void hash_init(int len,T prime=0xdefaced){
9     h_base[0]=1;
10    for(int i=1;i<=len;++i){
11        h[i]=(h[i-1]*prime+s[i-1])%prime_mod;
12        h_base[i]=(h_base[i-1]*prime)%prime_mod;
13    }
14 }
15 inline T get_hash(int l,int r){/*閉區間寫法，設編號為0 ~ len-1*/
16     return (h[r+1]-(h[l]*h_base[r-l+1])%prime_mod+prime_mod)%prime_mod;
17 }
```

Source: [日月卦長的模板庫](#)

→ [\[Rabin-Karp rolling hash \] Rabin-Karp 字串hash演算法](#)

Outline



KMP

- Knuth-Morris-Pratt algorithm
- 再來看個例子
a = “aabaac...”
b = “aabaab”

a	a	b	a	a	c	?	?	...
a	a	b	a	a	b			



KMP

- Knuth-Morris-Pratt algorithm
- 再來看個例子
a = “aabaac...”
b = “aabaab”

a	a	b	a	a	c	?	?	...
a	a	b	a	a	b			



KMP

- Knuth-Morris-Pratt algorithm
- 再來看個例子
a = “aabaac...”
b = “aabaab”

a	a	b	a	a	c	?	?	...
	a	a	b	a	a	b		



KMP

- Knuth-Morris-Pratt algorithm
- 再來看個例子
a = “aabaac...”
b = “aabaab”

a	a	b	a	a	c	?	?	...
		a	a	b	a	a	b	



KMP

- Knuth-Morris-Pratt algorithm

- 再來看個例子

a = “aabaac...”

b = “aabaab”

a	a	b	a	a	c
		a	a	b	a	a	b	

重複匹配
失敗



KMP

- Knuth-Morris-Pratt algorithm
- 再來看個例子
a = “aabaac...”
b = “aabaab”

a	a	b	a	a	c	?	?	...
a	a	b	a	a	b			



KMP

- Knuth-Morris-Pratt algorithm
- 再來看個例子
a = “aabaac...”
b = “aabaab”

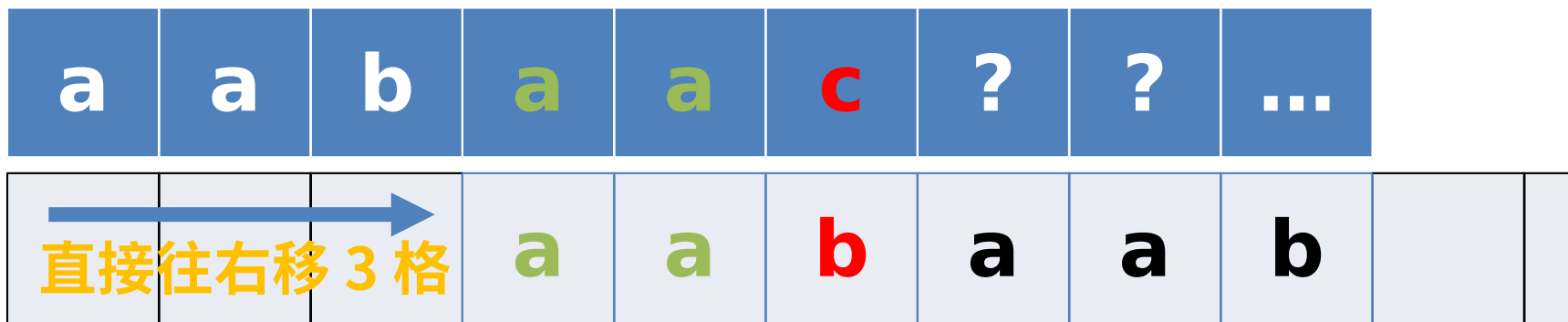


↑ ↑
早就可以知道他們一樣



KMP

- Knuth-Morris-Pratt algorithm
- 再來看個例子
a = “aabaac...”
b = “aabaab”



問題出在 B 有重複子字串

made by
Jingfei



KMP

- Knuth-Morris-Pratt algorithm

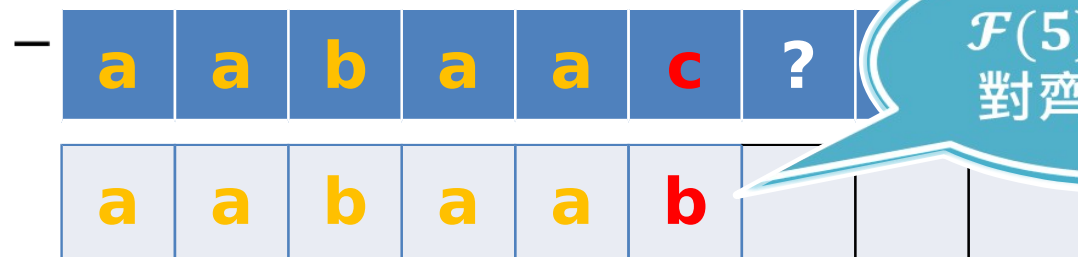
- 怎麼處理 B ?

- 定義 Fail function (失敗函數)

– 期望：能知道匹配失敗時，B 要對齊哪裡繼續匹配

$$\begin{aligned}
 & - \mathcal{F}_B(i) = \\
 & \quad \begin{cases} \max\{k: P_B(k) = B[0, k] = B[i - k, i]\}, & \text{if } i \neq 0 \text{ and at least a } k \text{ exists} \\ -1, & \text{else} \end{cases}
 \end{aligned}$$

– $\mathcal{F}(0) = -1$



$\mathcal{F}(5)$: 給我對齊B[2]!!

made by Jingfei



KMP

- Knuth-Morris-Pratt algorithm

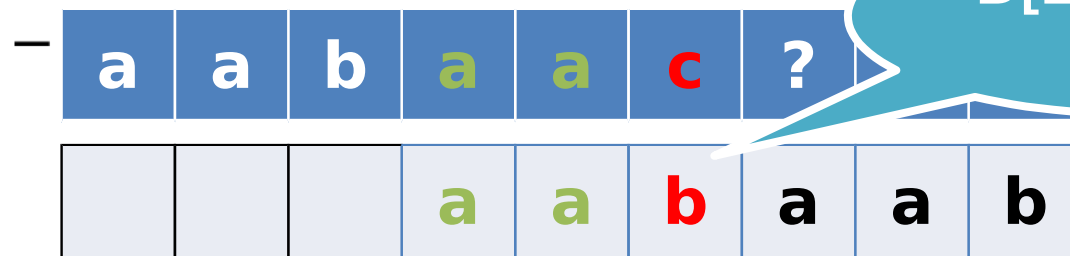
- 怎麼處理 B ?

- 定義 Fail function (失敗函數)

– 期望：能知道匹配失敗時，B 要對齊哪裡繼續匹配

$$F_B(i) = \begin{cases} \max\{k: P_B(k) = B[0, k] = B[i - k, i]\}, & \text{if } i \neq 0 \text{ and at least a } k \text{ exists} \\ -1, & \text{else} \end{cases}$$

– $F(0) = -1$



B[2] : 來了 ...



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1								

init:

$pi[0] = -1$

$cur_pos = -1$

made by
Jingfei



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1							

$B[\text{cur_pos}+1] \neq B[i]$
 $\text{pi}[i] = \text{cur_pos}$



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1						

$B[\text{cur_pos}+1] \neq B[i]$
 $\text{pi}[i] = \text{cur_pos}$



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1	0					

```

B[cur_pos+1]==B[i]
pi[i]=++cur_pos

```



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1	0	1				

```

B[cur_pos+1]==B[i]
pi[i]=++cur_pos

```



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1	0	1	2			

```
B[cur_pos+1]==B[i]
pi[i]=++cur_pos
```



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1	0	1	2	3		

```

B[cur_pos+1]==B[i]
pi[i]=++cur_pos

```



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1	0	1	2	3	4	

```

B[cur_pos+1]==B[i]
pi[i]=++cur_pos
    
```



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1	0	1	2	3	4	1

$B[\text{cur_pos}+1] \neq B[i]$
 $\text{pi}[i] = \text{pi}[\text{cur_pos}]$



KMP

- Fail function example

i	0	1	2	3	4	5	6	7	8
B	a	b	z	a	b	z	a	b	c
pi	-1	-1	-1	0	1	2	3	4	-1

$B[\text{cur_pos}+1] \neq B[i]$
 $\text{pi}[i] = \text{pi}[\text{cur_pos}]$



KMP

- Fail function

```
inline void fail (char *B , int *pi) {  
    int len = strlen (B );  
    pi[0] = -1;  
    for(int i = 1 , cur_pos = -1 ; i <= len ; ++ i) {  
        while (~ cur_pos && B [i] != B [cur_pos+ 1 ])  
            cur_pos = pi[cur_pos];  
        if(B [i] == B [cur_pos+ 1 ]) ++ cur_pos;  
        pi[i] = cur_pos;  
    }  
}
```



KMP

- Matching
- Fail function: 找出各後綴與前綴一樣的最大值
- 如果後綴 = 前綴 → 可直接位移



KMP

- Matching
- Fail function: 找出各後綴與前綴一樣的最大值
- 如果後綴 = 前綴 → 可直接位移



KMP

- Matching

A **x a b z a b z a b z a b c d**

cur_pos

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1

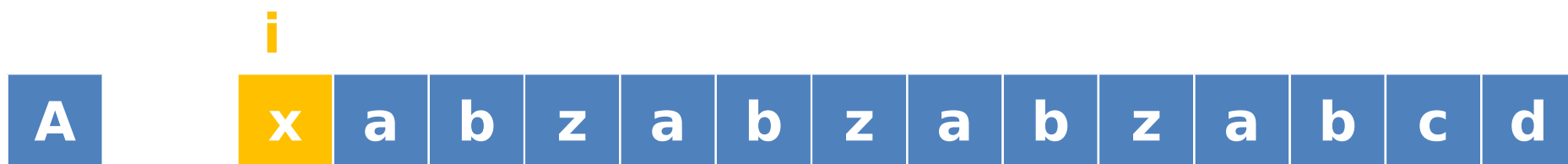
init:

cur_pos = -1



KMP

- Matching



cur_pos

$A[i] \neq B[\text{cur_pos} + 1]$

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1



KMP

- Matching



cur_pos

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1

$$A[i] == B[cur_pos + cur_pos]$$



KMP

- Matching



cur_pos

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	

$$A[i] == B[\text{cur_pos} + \text{cur_pos}]$$



KMP

- Matching



cur_pos

$A[i] == B[\text{cur_pos} + \text{cur_pos}]$

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	



KMP

- Matching



cur_pos

$A[i] == B[\text{cur_pos} + \text{cur_pos}]$

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	



KMP

- Matching



$A[i] == B[\text{cur_pos} + \text{cur_pos}]$

cur_pos

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	



KMP

- Matching



cur_pos

$A[i] == B[\text{cur_pos} + \text{cur_pos}]$

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	



KMP

- Matching



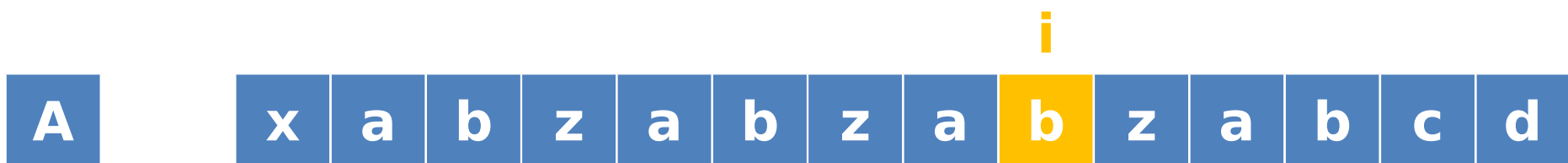
$A[i] == B[\text{cur_pos} + \text{cur_pos}]$

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	



KMP

- Matching



cur_pos

$A[i] == B[\text{cur_pos} + \text{cur_pos}]$

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1



KMP

- Matching

A x a b z a b z a b **z** a b c d

i

cur_pos

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1

$A[i] \neq B[\text{cur_pos} + 1]$
 $\text{cur_pos} = \text{pi}[\text{cur_pos}]$



KMP

- Matching



$$A[i] == B[\text{cur_pos} + \text{cur_pos}]$$

cur_pos

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	



KMP

- Matching



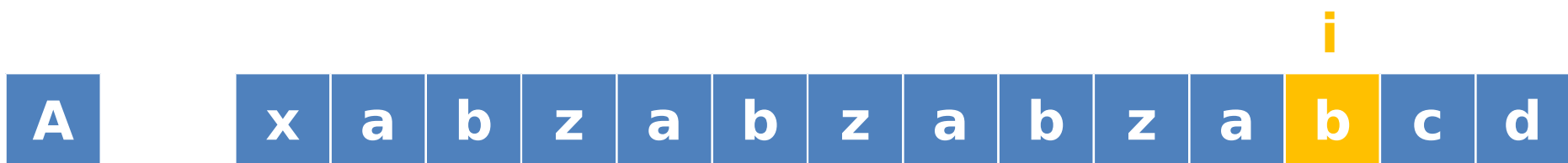
$$A[i] == B[cur_pos + cur_pos]$$

			cur_pos								
		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1



KMP

- Matching



$A[i] == B[\text{cur_pos} + \text{cur_pos}]$

cur_pos

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1



KMP

- Matching



cur_pos

$A[i] == B[cur_pos + cur_pos]$

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1



KMP

- Matching



		-1	0	1	2	3	4	5	6	7	8	
B			a	b	z	a	b	z	a	b	c	
pi			-1	-1	-1	0	1	2	3	4	-1	

cur_pos

A[i] == B[cur_pos + cur_pos]

cur_pos + 1 == len(B)

Match!!!



KMP

- Matching



cur_pos

		-1	0	1	2	3	4	5	6	7	8
B		a	b	z	a	b	z	a	b	c	
pi		-1	-1	-1	0	1	2	3	4	-1	

$A[i] == B[\text{cur_pos} + 1 + \text{cur_pos}]$
 $\text{cur_pos} + 1 == \text{len}(B)$
 $\text{cur_pos} = \text{pi}[\text{cur_pos}]$



KMP

- Matching



$A[i] \neq B[\text{cur_pos} + 1]$

cur_pos

		-1	0	1	2	3	4	5	6	7	8
B			a	b	z	a	b	z	a	b	c
pi			-1	-1	-1	0	1	2	3	4	-1



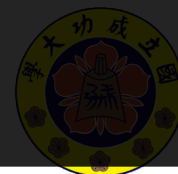
KMP

- Matching

```

inline void match(char *A, char *B, int *pi) {
    int lenA = strlen(A);
    int lenB = strlen(B);
    for(int i = 1, cur_pos = -1; i <= lenA; ++ i) {
        while (~ cur_pos && A[i] != B[cur_pos + 1])
            cur_pos = pi[cur_pos];
        if(A[i] == B[cur_pos + 1]) ++ cur_pos;
        if(cur_pos + 1 == lenB) {
            /* Match !! */
            cur_pos = pi[cur_pos];
        }
    }
}

```



KMP

- Fail function + Matching
- Complexity
 - 關鍵: while-loop
 - cur_pos 每次只會 +1 或往前
 - 均攤後 $O(|A| + |B|)$

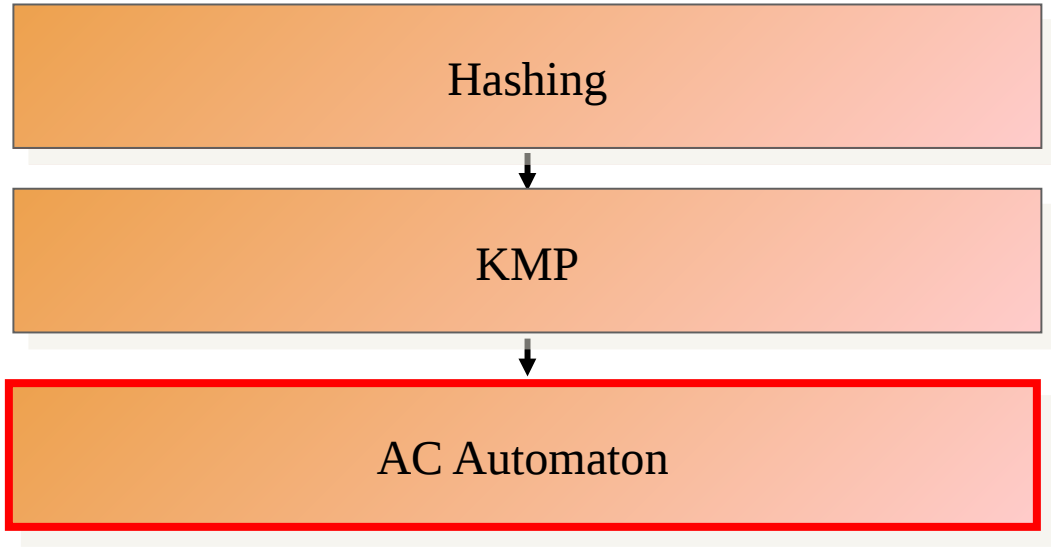


Example

- [POJ 3461](#)
- [UVA 455](#)



Outline

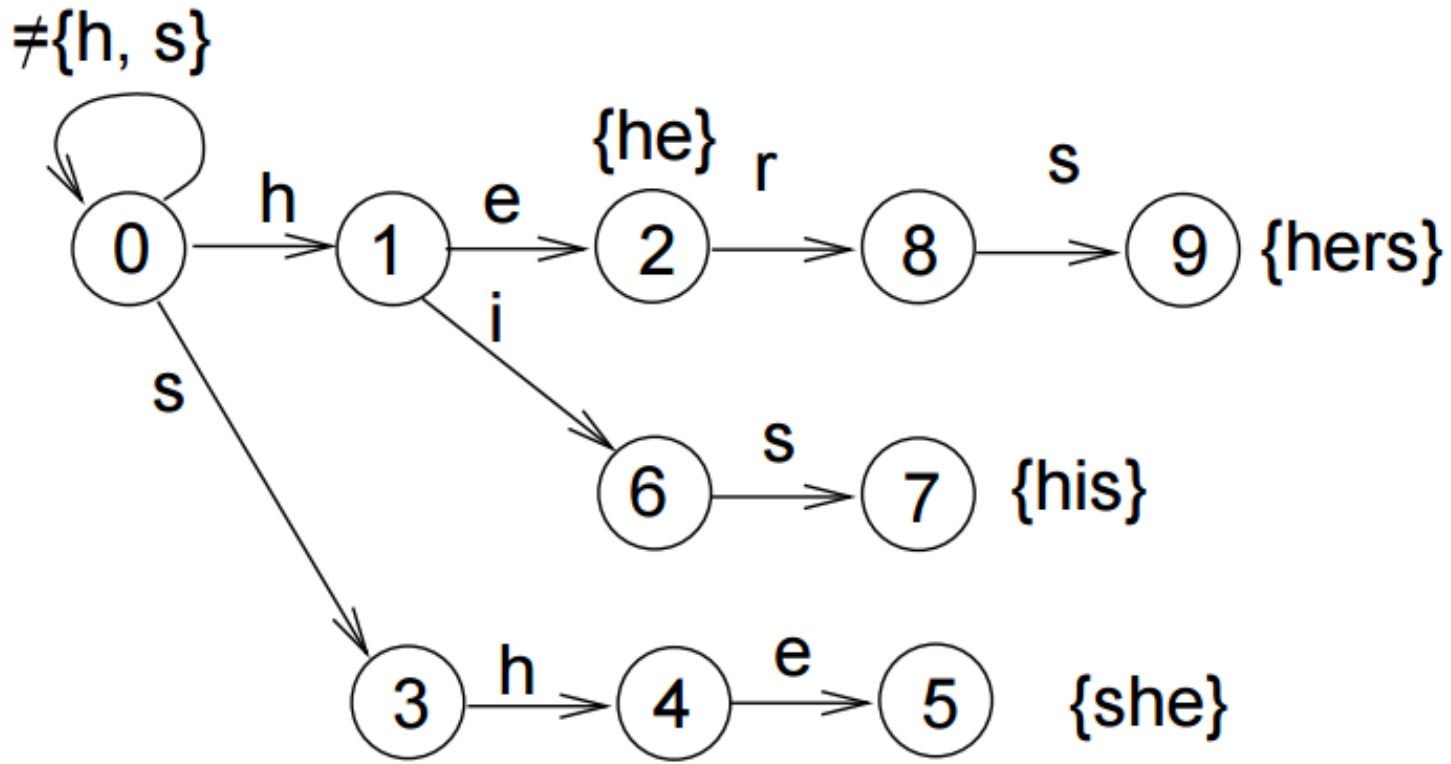


AC Automaton

- KMP 複雜度 $O(|A| + |B|)$
- 多字串匹配
 1. 一個字串 B 匹配很多字串 A_i
 - $\Rightarrow O(\sum |A_i| + |B|)$
 - \Rightarrow 線性
 2. 很多字串 B_i 匹配一個字串 A
 - $\Rightarrow O(n|A| + \sum |B_i|)$
 - \Rightarrow 弱弱的
- Trie : 儲存多個字串
- ACC 自動機 \equiv KMP + Trie



AC Automaton



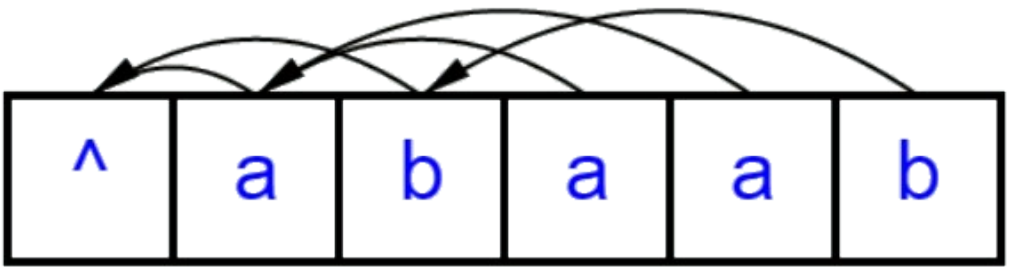
- Trie : 儲存多個字串
- ACC 自動機 ≡ KMP + Trie



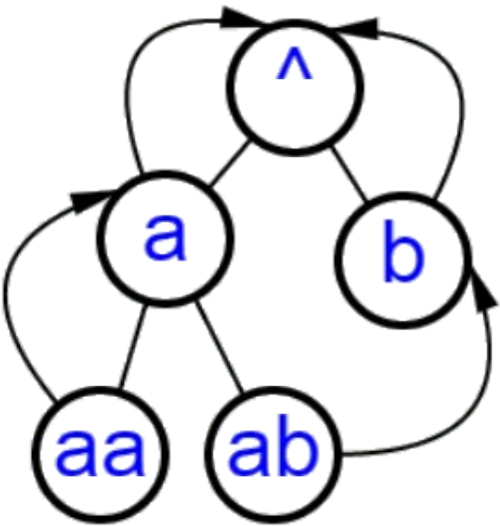
AC Automaton

- 比較 Fail function (圖)

KMP:



AC自動機:



AC Automaton

- 比較 Failfunction (定義)

- KMP

- $\mathcal{F}_B(i) =$

- $\max\{k: B[0, k] = B[i - k, i]\}, \text{ if } i \neq 0 \text{ and at least a } k \text{ exists}$
 - $-1, \text{ else}$

- AC Automaton

- $A[0, k]$ 是 $A[0, i]$ 的前綴

- AC Automaton

- $\mathcal{F}_B(v) = \begin{cases} u, & \text{if } B_T(u) \text{ 是 } B_T(v) \text{ 的前綴且 } |S_T(u)| \text{ 最大} \\ v_0, & \text{else} \end{cases}$

- $B_T(u)$ 是 $B_T(v)$ 的前綴



AC Automaton

- 比較 Fail function (匹配失敗)
- KMP
 - 沿著 $\mathcal{F}(i)$, 嘗試 (i) , 直到嘗試, 直到 $\mathcal{F}^t(i) = -1$
- AC Automaton
 - 沿著 $\mathcal{F}(v)$, 嘗試 (v) , 直到嘗試, 直到 $\mathcal{F}^t(v) = v_0$ (v_0 : root)



AC Automaton

- 比較 Fail function (構造)
- KMP
 - 利用 $\pi(i-1)$ 求出 $\mathcal{F}(i)$
- AC Automaton
 - 利用 $\pi(u)$ 求出 $\pi(v)$ 的父節點
 - use BFS



Automaton – 範例 code

- 田用卦長的模板庫

 - ↳ Aho Corasick Automaton | AC自動機

- AC自动机算法详解

- AC自动机基础入门(PPT) by 李翔
- AC自动机基础入门(PPT) by 李翔



Automaton – 例題

- AC Automaton
 - [SPOJ NSUBSTR](#) 題解
 - [SPOJ SUBLEX](#) 題解
 - [Codeforces 235C](#) 題解
 - 這些都是作法很多 (Suffix Array, Suffix Tree...) ，非常經典的問題，可以從中理解 SAM 的精妙之處
- UVa 10679, 1449



Reference

- 歷屆 PPT…… (electron, free999, louis6340, …)
- 2015 IOI camp 字串處理
<http://ioicamp.csie.org/content>
<http://bobogei81123.github.io/ioi-lecture>

