

NCKU Programming Contest Training Course

2016/04/20

Jingfei Yang

<http://myweb.ncku.edu.tw/~e84016184/StringMatching.pdf>

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 $A = a_0a_1\dots a_{n-1}$
 - $a_i \in$ 字元集， n 是字串的長度
- 子字串 substring
 - $A[i, j] = a_ia_{i+1}a_{i+2} \dots a_j$ (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 $S_A(k) = a_ka_{k+1}a_{k+2} \dots a_n, 0 \leq k < n$
- 前綴 prefix
 - A 的一個子字串 $P_A(h) = a_0a_1a_2 \dots a_h, 0 \leq h < n$



String Basic

- $S = \text{"abcbbab"}$
 - 子字串 : "bcb" , "bba" , ...
 - 子序列 : "acb" , "bbb" , ...
 - 前綴 : "abcb" , "ab" , ...
 - 後綴 : "bbab" , "ab" , ...



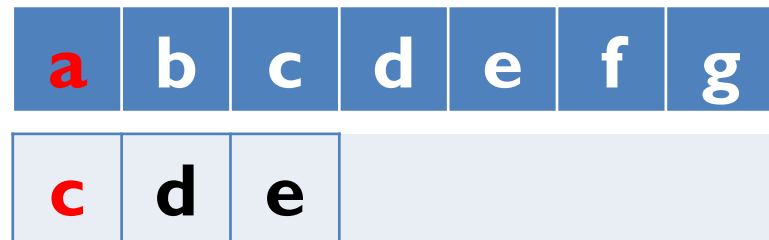
String Matching

~	!	@	#	\$	%	^	&	*	()	-	+	=	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： $\mathcal{O}(|A|)$
- A = "abcdefg"
B = "cde"



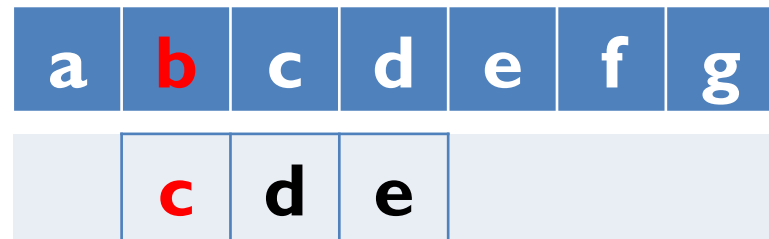
String Matching

~	!	@	#	\$	%	^	&	*	()	-	=	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： $\mathcal{O}(|A|)$
- A = "abcdefg"
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： $\mathcal{O}(|A|)$
- A = "abcdefg"
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： $\mathcal{O}(|A|)$
- A = "abcdefg"
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 }
```

Matching!!

- 複雜度： $\mathcal{O}(|A|)$
- A = "abcdefg"
B = "cde"



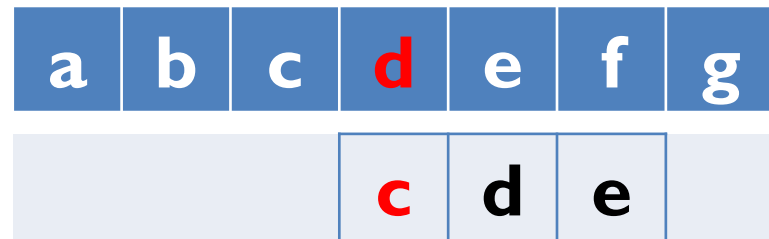
String Matching

~	!	@	#	\$	%	^	&	*	()	-	+	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： $\mathcal{O}(|A|)$
- A = "abcdefg"
B = "cde"



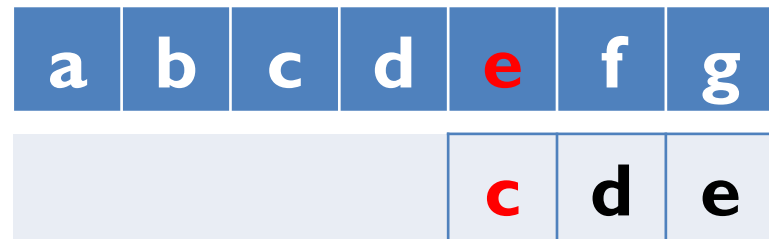
String Matching

~	!	@	#	\$	%	^	&	*	()	-	+	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： $\mathcal{O}(|A|)$
- A = "abcdefg"
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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"

a	a	a	a	a	a
a	a	a	b		



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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"

a	a	a	a	a	a
a	a	a	b		



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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"

a	a	a	a	a	a
a	a	a	b		



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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"

a	a	a	a	a	a
a	a	a	b		



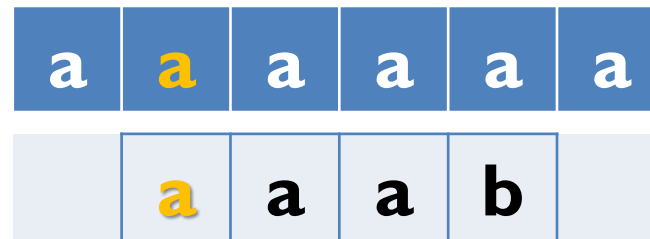
String Matching

~	!	@	#	\$	%	^	&	*	()	_	+	=	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： $\mathcal{O}(|A|)$
- 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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"



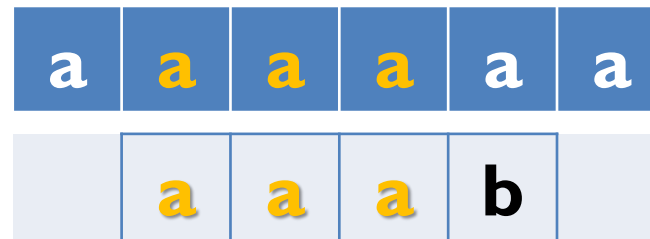
String Matching

~	!	@	#	\$	%	^	&	*	()	_	+	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： $\mathcal{O}(|A|)$
- 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： $\mathcal{O}(|A|)$
- 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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"

a	a	a	a	a	a
		a	a	a	b



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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"

a	a	a	a	a	a
		a	a	a	b



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： $\mathcal{O}(|A|)$
- A = "aaaaaaaa...aaa"
B = "aaaaaaaa...aab"

a	a	a	a	a	a
		a	a	a	b



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 }
```

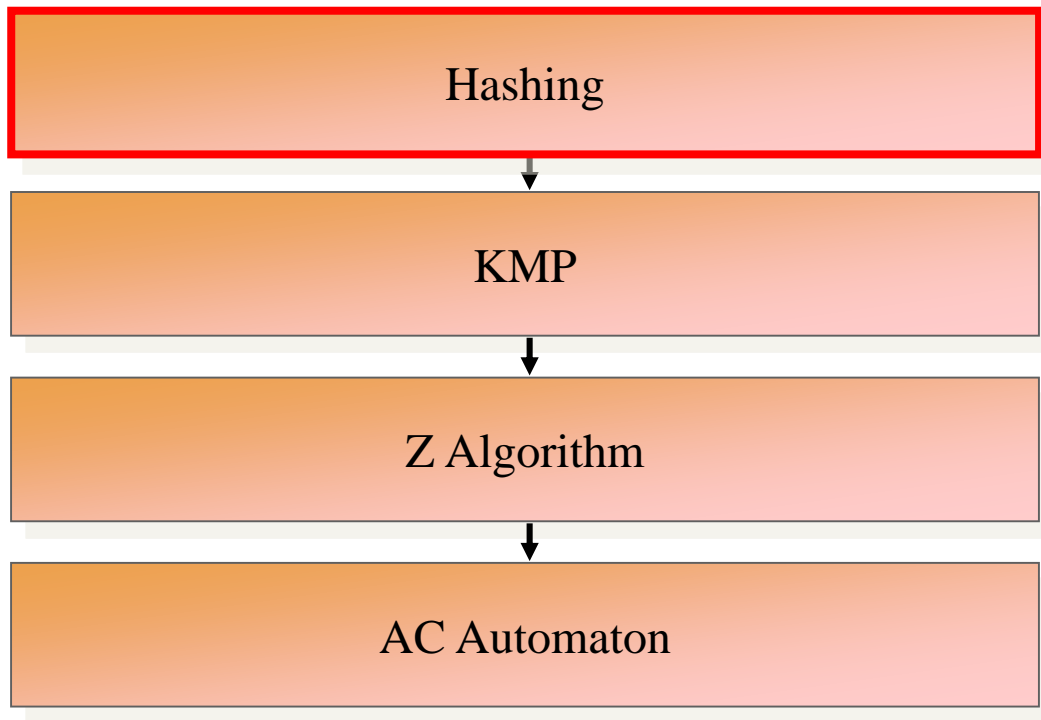
- 複雜度： ~~$O(|A|)$~~ $O(|A||B|)$

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

a	a	a	a	a	a
		a	a	a	b



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
 - f 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

- $f(A) = a_0p^{n-1} + a_1p^{n-2} + \dots + a_{n-2}p + a_{n-1} \bmod q$
- 類似： p 進位制，分成 q 類
- p, q 取不同質數 \rightarrow 均勻

- 滾動

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

\rightarrow 計算 A 所有前綴的 hash value， $O(|A|)$

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

\rightarrow 任何 A 子字串的 hash value， $O(1)$

3. 枚舉 A 長度為 $|B|$ 的子字串，比較 hash value

$\rightarrow O(N)$



Hashing

- A = "abcdefg"

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

- B = "cde" = $2p^2 + 3p^1 + 4p^0$



Hashing

- A = "abcdefg"

a	b	c	d	e	f	g
	$1p^2 + 2p^1 + 3p^0$					

- B = "cde" = $2p^2 + 3p^1 + 4p^0$



Hashing

- A = "abcdefg"

a	b	c	d	e	f	g
		$2p^2 + 3p^1 + 4p^0$				

- B = "cde" = $2p^2 + 3p^1 + 4p^0$

Matching!!



Hashing

- A = "abcdefg"

a	b	c	d	e	f	g
			$3p^2 + 4p^1 + 5p^0$			

- B = "cde" = $2p^2 + 3p^1 + 4p^0$



Hashing

- A = "abcdefg"

a	b	c	d	e	f	g
				$4p^2 + 5p^1 + 6p^0$		

- B = "cde" = $2p^2 + 3p^1 + 4p^0$

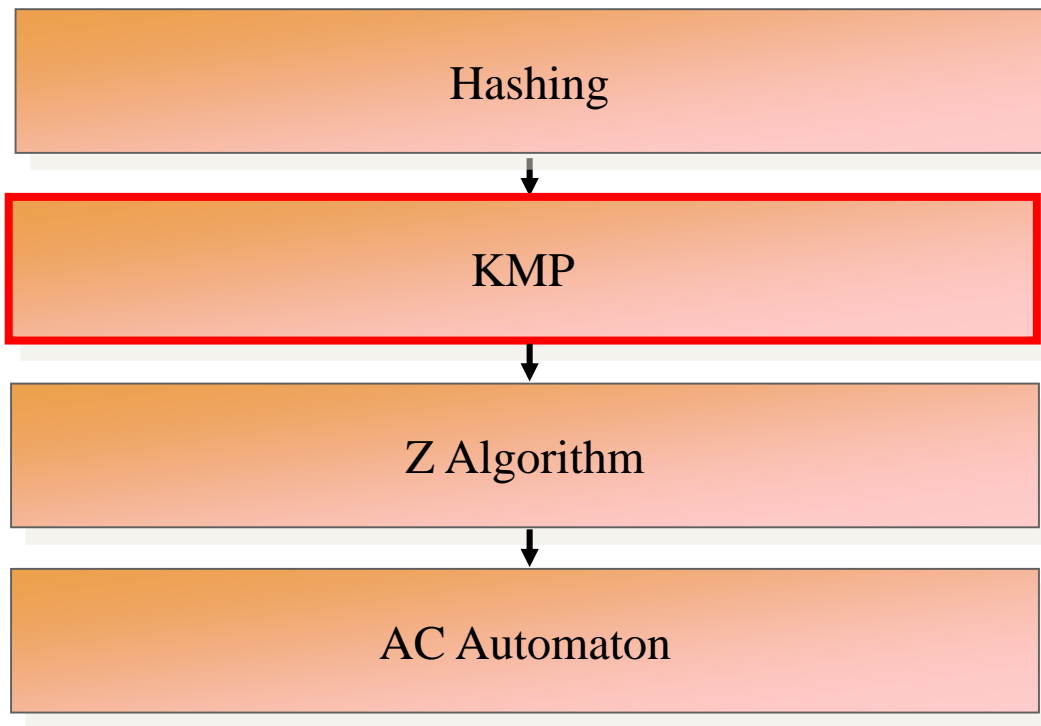


Hashing

- 回到剛剛
 - $A \neq B \Rightarrow f(A) \neq f(B)$ → 不一定
 - 相等時重新檢查一次？
 - $A = \text{"aaaaaaaa...aaa"}$
 $B = \text{"aaaaaaaa...aab"}$
- q 取大一點 (long long 質數)
 - 碰撞機率小
- ex. $q \in 10^{15} \Rightarrow \text{probability: } 10^{-15}$
- ex. 2147483647



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"

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

直接往右移3格

a	a	b	a	a	b
---	---	---	---	---	---

↑ ↑
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"

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


直接往右移3格 →			a	a	b	a	a	b	
-----------	--	--	---	---	---	---	---	---	--

問題出在 B 有重複子字串



KMP

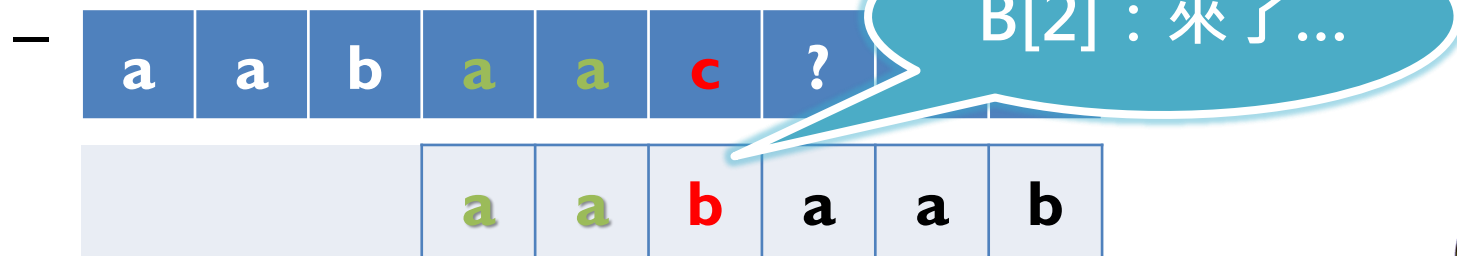
- Knuth-Morris-Pratt algorithm
- 怎麼處理 B ?
- 定義 **Fail function** (失敗函數)
 - 期望： $\mathcal{F}(i)$ 能知道匹配失敗時，B 要對齊哪裡繼續匹配
 - $\mathcal{F}_B(i) = \begin{cases} \max\{k: P_B(k) = B[0, k] = B[i - k, i]\}, & \text{if } i \neq 0 \text{ and at least a } k \text{ exists} \\ -1, & \text{else} \end{cases}$
 - $\mathcal{F}(0) = -1$
 - | | | | | | | |
|---|---|---|---|---|---|---|
| a | a | b | a | a | c | ? |
| a | a | b | a | a | b | |



$\mathcal{F}(5)$: 給我對齊B[2]!!

KMP

- Knuth-Morris-Pratt algorithm
- 怎麼處理 B ?
- 定義 **Fail function** (失敗函數)
 - 期望： $\mathcal{F}(i)$ 能知道匹配失敗時，B 要對齊哪裡繼續匹配
 - $\mathcal{F}_B(i) = \begin{cases} \max\{k: P_B(k) = B[0, k] = B[i - k, i]\}, & \text{if } i \neq 0 \text{ and at least a } k \text{ exists} \\ -1, & \text{else} \end{cases}$
 - $\mathcal{F}(0) = -1$



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					

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



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	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			

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



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	1	2	3		

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



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	1	2	3	4	

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



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	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

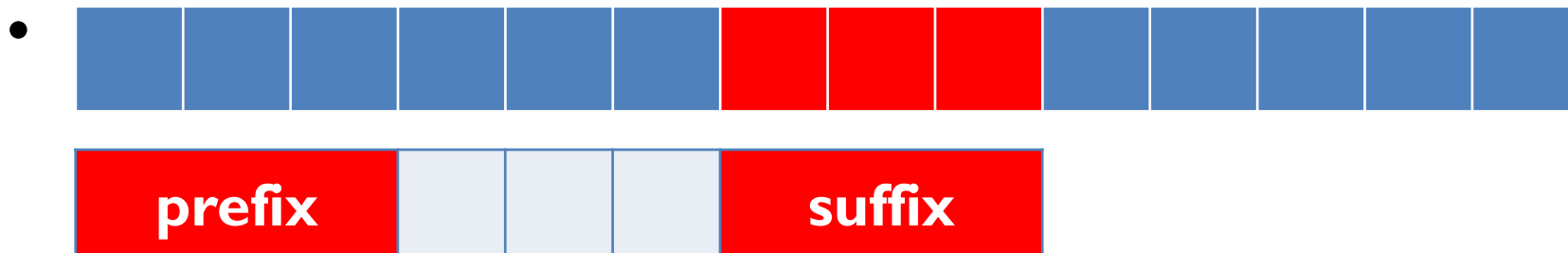
- Fail function

```
1 void fail(string B, int *pi){
2     int len = strlen(B);
3     pi[0] = -1;
4     for(int i=1, cur_pos=-1; i<len; ++i){
5         while(cur_pos>=0 && B[i]!=B[cur_pos+1])
6             cur_pos=pi[cur_pos];
7         if(B[i]==B[cur_pos+1]) ++cur_pos;
8         pi[i]=cur_pos;
9     }
10 }
```



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

A ⁱ
x a b z a b z a b z a b c d

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

i

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

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



KMP

- Matching

			i											
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

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



KMP

- Matching

A

x

a

b

z

a

b

z

a

b

z

a

b

c

d

i

B

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

cur_pos

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



KMP

- Matching

A

x

a

b

z

a

b

z

a

b

z

a

b

c

d

i

B

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

cur_pos

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



KMP

- Matching

i

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

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



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

a

b

z

a

b

z

a

b

c

-1

-1

-1

0

1

2

3

4

-1

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

++cur_pos

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



KMP

- Matching

i

A	x	a	b	z	a	b	z	a	b	z	a	b	c	d
----------	---	---	---	---	---	---	---	---	---	---	---	---	---	---

	-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

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



KMP

- Matching

A

x

a

b

z

a

b

z

a

b

z

a

b

c

d

i

B

pi

-1

0

1

2

3

4

5

6

7

8

a

b

z

a

b

z

a

b

c

-1

-1

-1

0

1

2

3

4

-1

cur_pos

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

++cur_pos

$A[i] == B[cur_pos + 1]$
 $++cur_pos$



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

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

$$A[i] == B[cur_pos + 1]$$

$$++cur_pos$$



KMP

- Matching

A	x	a	b	z	a	b	z	a	b	z	a	b	c	d
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

i

	-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+1]
++cur_pos

- A horizontal array of 14 blue boxes containing the characters 'x', 'a', 'b', 'z', 'a', 'b', 'z', 'a', 'b', 'z', 'a', 'b', 'c', 'd'. Above the 12th box (containing 'b'), there is a yellow box containing the letter 'i', representing the current index.

	-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+1]
++cur_pos
```

KMP

- Matching

														<i>i</i>	
A	x	a	b	z	a	b	z	a	b	z	a	b	c	d	

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



KMP

- Matching

A

x

a

b

z

a

b

z

a

b

z

a

b

c

d

i

	-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+1]

++cur_pos

cur_pos+1==len(B)

Match!!!



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

$A[i] == B[cur_pos + 1]$
 $++cur_pos$
 $cur_pos + 1 == len(B)$
 $cur_pos = pi[cur_pos]$



KMP

- Matching

A x a b z a b z a b z a b c **d** ⁱ

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

```
1 void fail(string A, string B, int *pi){
2     int lenA = A.length();
3     int lenB = B.length();
4     for(int i=0, cur_pos=-1; i<lenA; ++i){
5         while(cur_pos>=0 && A[i]!=B[cur_pos+1])
6             cur_pos=pi[cur_pos];
7         if(A[i]==B[cur_pos+1]) ++cur_pos;
8         if(cur_pos+1==lenB){
9             /* Match!!! */
10            cur_pos=pi[cur_pos];
11        }
12    }
13 }
```

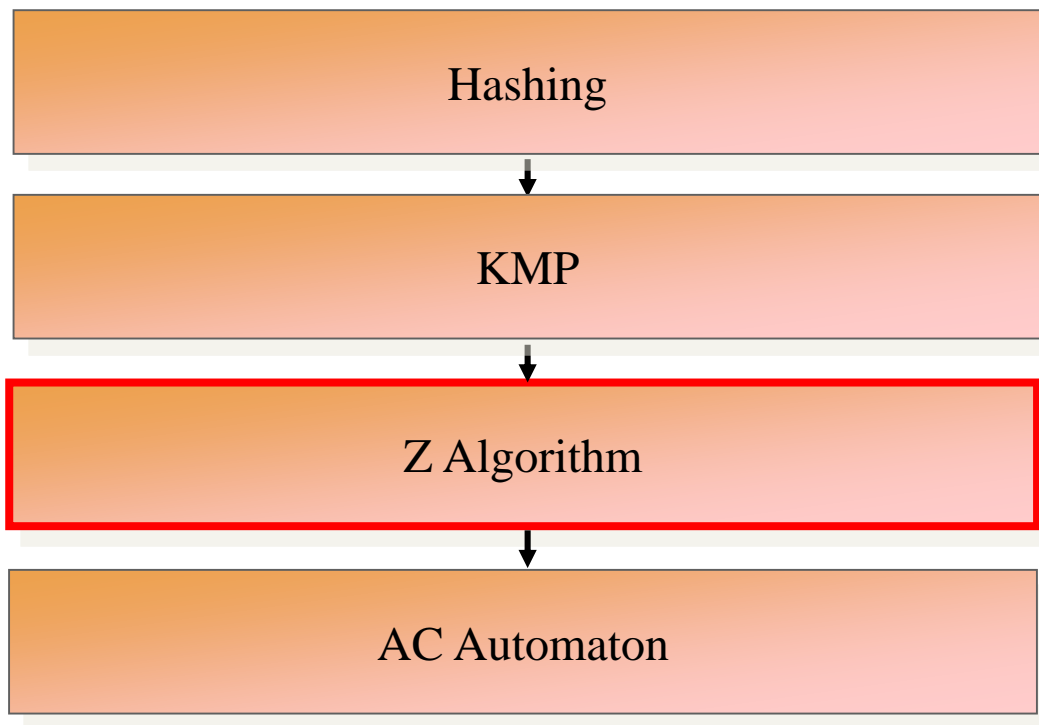


KMP

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

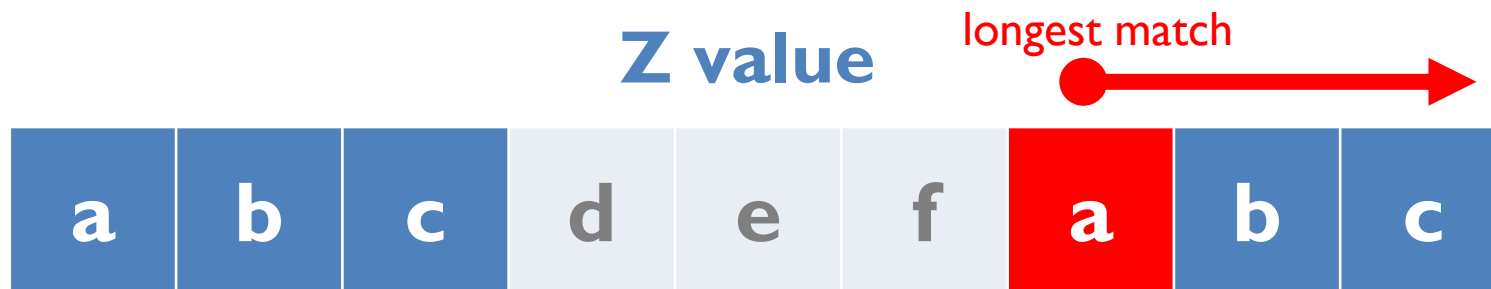
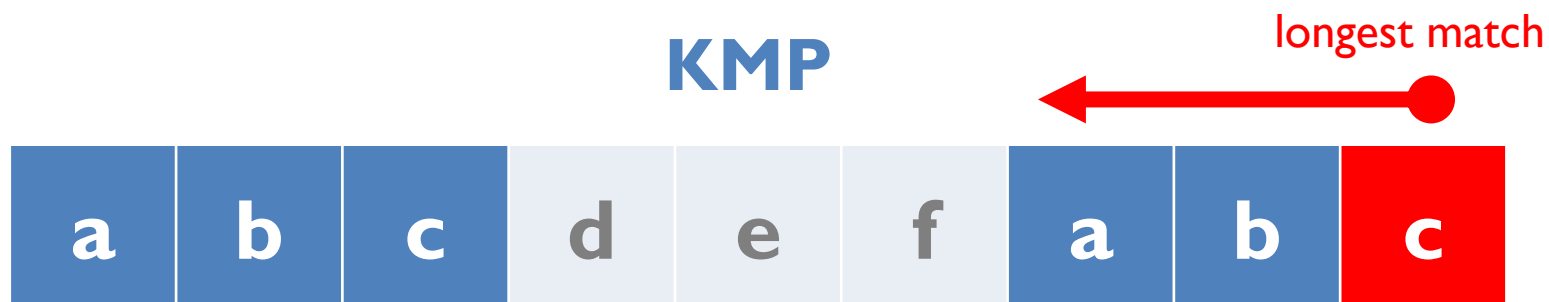


Outline



Z Algorithm

- vs KMP
 - S= "abcdefabc"



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z									



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9								



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0							



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0	0						



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0	0	1					



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0	0	1	4				



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0	0	1	4	0			



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0	0	1	4	0	0		



Z Algorithm

- 定義

- $Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0	0	1	4	0	0	2	



Z Algorithm

- 定義

- $$Z_A(i) = \begin{cases} 0, & \text{if } i = 0 \text{ or } A[i] \neq A[0] \\ \max\{k: A[0, k-1] = A[i, i+k-1]\}, & \text{else} \end{cases}$$
- 由 $A[i]$ 開始的字串，可以和 A 自己匹配多長
- ex.

i	0	1	2	3	4	5	6	7	8
A	a	b	z	a	a	b	z	a	b
z	9	0	0	1	4	0	0	2	0



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$
 - $A[z] \neq A[i + z]$

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$
 - $A[z] \neq A[i + z]$
 - 令 $L = i$, $R = i + z - 1$, $L \leq j \leq R$, $j' = j - L$

							L						R	
i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



made by Jingfei



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$
 - $A[z] \neq A[i + z]$
 - 令 $L = i$, $R = i + z - 1$, $L \leq j \leq R$, $j' = j - L$
 - case $\times 3$

	<i>L</i>							<i>R</i>						
i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



made by Jingfei



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$
 - $A[z] \neq A[i + z]$
 - 令 $L = i$, $R = i + z - 1$, $L \leq j \leq R$, $j' = j - L$
 - **case 1.** $j' + Z(j') < z \Rightarrow Z(j) = Z(j')$

							L						R	
i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



made by Jingfei



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$
 - $A[z] \neq A[i + z]$
 - 令 $L = i$, $R = i + z - 1$, $L \leq j \leq R$, $j' = j - L$
 - **case 2.** $j' + Z(j') > z \Rightarrow Z(j) = R - j + 1$ (j 到 R 的長度)

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



made by Jingfei



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$
 - $A[z] \neq A[i + z]$
 - 令 $L = i$, $R = i + z - 1$, $L \leq j \leq R$, $j' = j - L$
 - **case 3.** $j' + Z(j') = z \Rightarrow Z(j) \geq Z(j')$ (剛好在邊界)

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



made by Jingfei



Z Algorithm

- 假設一個字串 A 中， $Z(i) = z$
 - $A[k] = A[i + k]$, if $0 \leq k < z$
 - $A[z] \neq A[i + z]$
 - **更新** $L = j$, $R = j + Z(j) - 1$, $L \leq j \leq R$, $j' = j - L$
 - **case 3.** $j' + Z(j') = z \Rightarrow Z(j) \geq Z(j')$

	L										R			
i	0	1	2	3	4	5	6	7	8	9	10	11	12	13
A	a	b	a	b	a	b	e	a	b	a	b	a	b	f
z	14	0	4	0	2	0	0	6	0	4	0	2	0	0



made by Jingfei



Z Algorithm

- 只會前進不會後退 $\rightarrow O(N)$

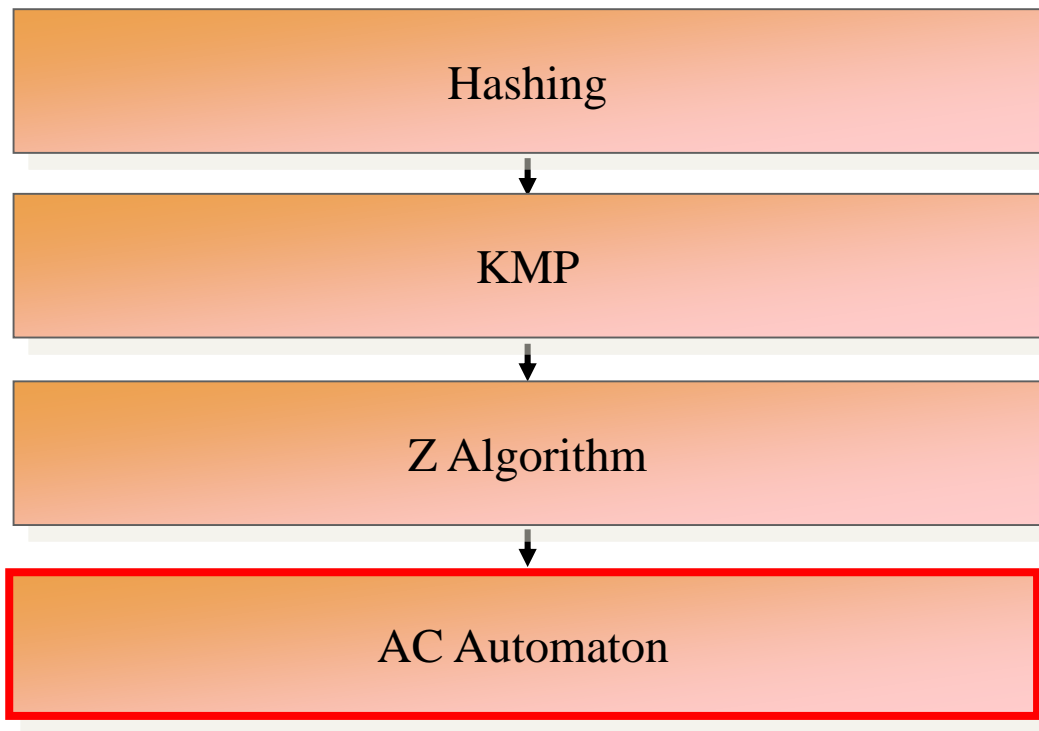
```
1 L = R = 0;
2 for(int i=1; i<len; i++){
3     if(i>R) Z[i]=0; //Case 0
4     else{
5         int ip = i - L;
6         if(ip+Z[ip] < Z[L]) Z[i]=Z[ip]; //Case 1
7         else Z[i]=R-i+1; //Case 2, 3
8     }
9     while(i+Z[i] < len && A[i+Z[i]]==A[Z[i]])
10         Z[i]++;
11     if(i+Z[i]-1 > R){
12         L=i;
13         R=i+Z[i]-1;
14     }
15 }
```

Z Algorithm

- Matching ?
- 把兩個字串接起來， B 在前， A 在後，
中間用一個沒有出現過的字元連接
- 如 $A = \text{"abaab"}, B = \text{"aab"}$
令 $S = \text{"aab$abaab"}$
→ 發現 B 是否能在 A 的某個位置被匹配，
只要看那個位置的 Z value 是否等於 B 的長度
- 最長回文子字串



Outline

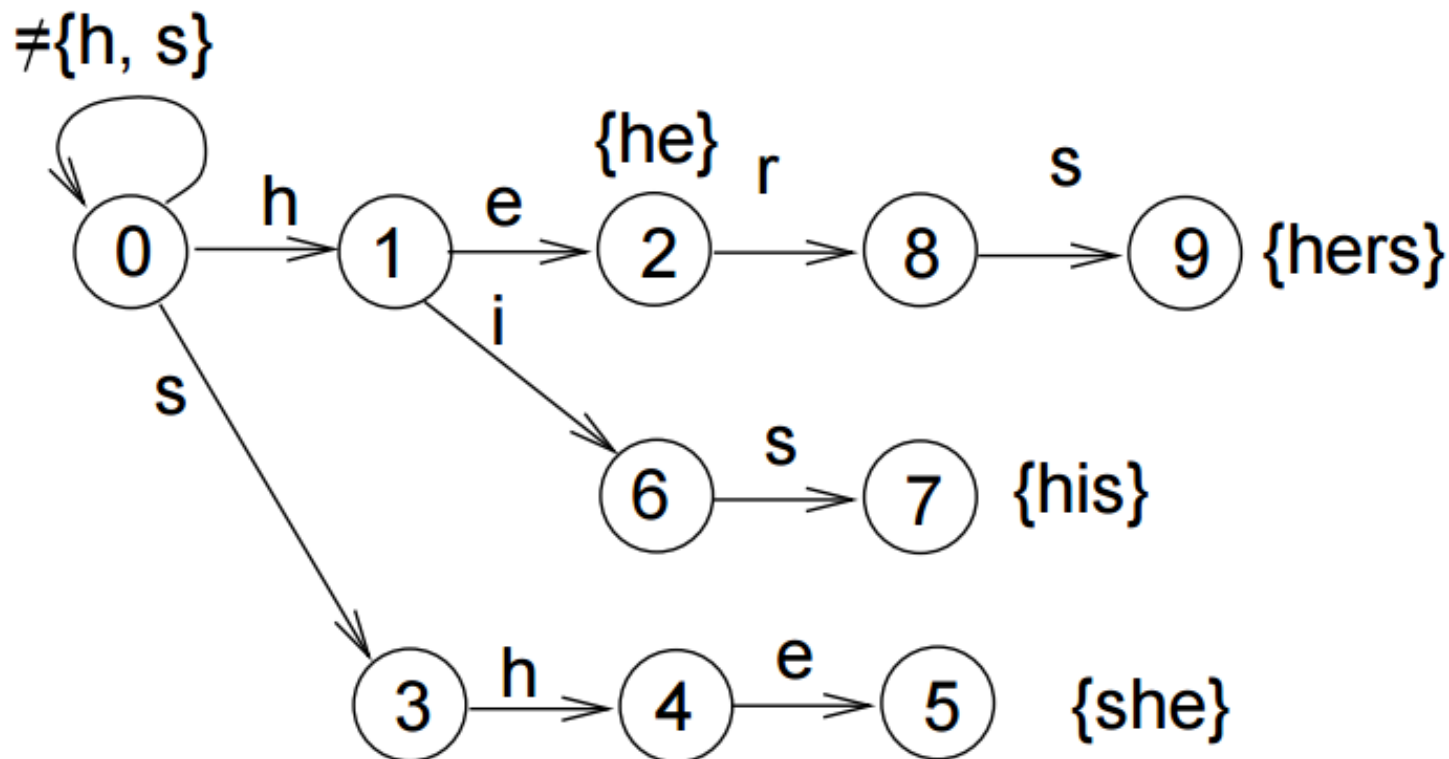


AC Automaton

- 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
 - $O(\sum |A_i| + |B|)$
 - 線性
 2. 很多字串 B_i 匹配一個字串 A
 - $O(n|A| + \sum |B_i|)$
 - 弱弱的
- Trie : 儲存多個字串
- AC 自動機 = KMP + Trie



AC Automaton



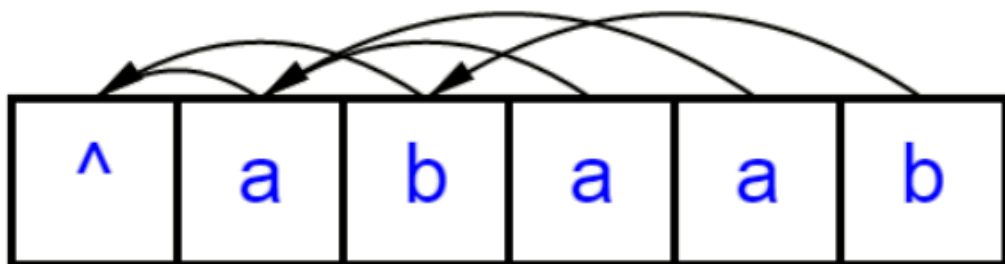
- Trie : 儲存多個字串
- AC 自動機 = KMP + Trie



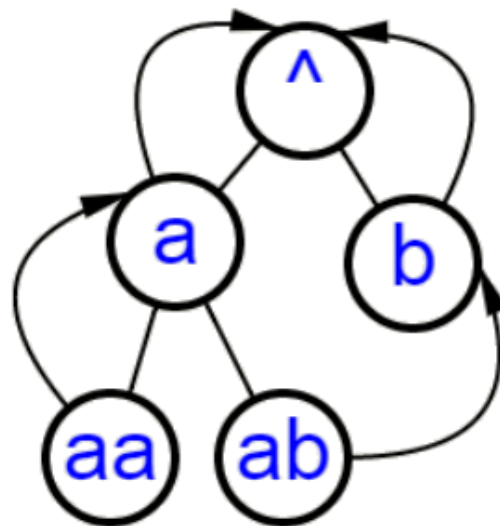
AC Automaton

- 比較 Fail function (圖)

KMP:



AC自動機:



AC Automaton

- 比較 Fail function (定義)
- KMP
 - $\mathcal{F}_B(i) = \begin{cases} \max\{k: P_B(k) = B[0, k] = B[i - k, i]\}, & \text{if } i \neq 0 \text{ and at least a } k \text{ exists} \\ -1, & \text{else} \end{cases}$
 - $A[0, k]$ 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 $\mathcal{F}(i)$, $\mathcal{F}^2(i)$, ... 嘗試，直到 $\mathcal{F}^t(i) = -1$
- AC Automaton
 - 沿著 $\mathcal{F}(v)$, $\mathcal{F}^2(v)$, ... 嘗試，直到 $\mathcal{F}^t(v) = v_0$ (v_0 : root)



AC Automaton

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



AC Automaton

- $O(|A| + \sum |B|)$

```
1 root->fail = NULL;
2 queue< Node* > que;
3 que.push_back(root);
4 while ( !que.empty() ) {
5     Node *fa = que.front(); que.pop_front();
6
7     for (auto it = fa->child.begin();
8         it != fa->child.end(); it++) {
9         Node *cur = it->second, *ptr = fa->fail;
10        while ( ptr && !ptr->child.count(it->first) )
11            ptr = ptr->fail;
12
13        cur->fail = ptr ? ptr->child[it->first] : root;
14        que.push(cur);
15    }
16 }
```



Example

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



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>

