

## 第三章 词法分析

### 本章内容

- 有限自动机
- 词法分析器的设计与实现
- 词法分析器的自动生成

2

### 3.1 有限自动机

- 确定有限自动机
- 非确定有限自动机
- 正规文法与确定自动有限自动机的等价性
- 正规式与确定自动有限自动机的等价性
- 确定自动有限自动机的化简

3

#### 3.1.1 正规式与正规集

- Regular expressions represents patterns of strings of characters
- The language generated by the regular expression  $r$  is the set of strings that it matches, and is written as  $L(r)$
- 字母表  $\Sigma$ ,
- 元字符 (元符号), metacharacters, metasymbols
- 托字符, escape character

4

## 基本正规式

- 任何  $a \in \Sigma$ ,  $a$  是  $\Sigma$  上的一个正规式, 它所表示的正规集为  $L(a) = \{a\}$ ;
- $\epsilon$  和  $\phi$  都是字母表  $\Sigma$  上的正规式, 它们所表示的正规集分别为  $\{\epsilon\}$  和  $\phi$ ;

5

## 正规式的运算

- $r$  和  $s$  是  $\Sigma$  上的正规式, 那么:
  - $\oplus$  选择运算 |  
 $r | s$  是正规式且  $L(r | s) = L(r) \cup L(s)$ ;
  - $\oplus$  连接运算 •  
 $r * s$  是正规式且  $L(r * s) = L(r) * L(s)$ ;
  - $\oplus$  重复运算 \*  
 $r^*$  是正规式且  $L(r^*) = L(r)^*$ ;
- 运算的优先级和括号使用

通过有限次使用上述三个步骤而得到的表达式是  $\Sigma$  上的正规式. 由这些正规式所表示的符号串的集合是  $\Sigma$  上的正规集.

6

## 运算的优先级

- $a | b^*$ 
    - $\oplus (a | b)^*$
    - $\oplus a | (b^*)$
  - 按惯例 \* 的优先级较高
  - $a | bc^*$ 
    - $\oplus a | (b(c^*))$
  - $ab | c^*d$ 
    - $\oplus (ab) | ((c^*)d)$
  - 按惯例 \* 的优先级最高, | 最低
- 可用括号改变运算次序  
 $\oplus (a | b)c$   
 $\oplus (a | bb)^*$

7

## 举例

- 已知  $\Sigma = \{a, b\}$ , 那么  $ba^*$ ,  $a(a | b)^*$  和  $(a | b)^*(aa | bb)(a | b)^*$  都是  $\Sigma$  上的正规式.
- $a$  是 (2, p46)  $\Rightarrow a^*$  是 (3 闭包),  $b$  是 (2)  $\Rightarrow ba^*$  是 (3 连接)  
 正规集:  $\Sigma$  上所有以  $b$  为前后跟 0 到多个  $a$  的符号串.
- $a$  是 (2),  $b$  是 (2)  $\Rightarrow a | b$  是 (3 或)  $\Rightarrow (a | b)^*$  (3 闭包) 是;  $a$  是 (2)  $\Rightarrow a(a | b)^*$  是 (3 连接)  
 正规集:  $\Sigma$  上所有以  $a$  为首的符号串.
- $a$  是 (2),  $b$  是 (2)  $\Rightarrow a | b$  是 (3 或)  $\Rightarrow (a | b)^*(3 \text{ 闭包})$  是;  $a$  是 (2)  $\Rightarrow aa$  是 (3 连接);  $b$  是 (2)  $\Rightarrow bb$  是 (3 连接)  $\Rightarrow (a | b)^*(aa | bb)(a | b)^*$  是 (3 连接)  
 正规集:  $\Sigma$  上所有含有两个相继  $a$  或两个相继  $b$  的符号串.

8

## 举例

令字母表 $\Sigma=\{a,b\}$ , 则:

正规式 $a|b$ 表示集合{a, b}

正规式 $(a|b)(a|b)$ 表示集合{aa, ab, ba, bb}

正规式 $a^+$ 表示由1到多个a构成的符号串的集合

正规式 $(a|b)^*$ 表示由0到多个a或b构成的符号串的集合

正规式 $a|a^*b$ 表示符号串a以及由0到多个a后跟一个b构成的符号串的集合

9

## 举例

对于定义在字母表{a,b}上的语言, 分别给出下述正规集的正规式:

1) 所有以a为开头和结尾的符号串  $\epsilon|a|a(a|b)^*a$

2) 第奇数位都是a的所有符号串  $(a(a|b))^*(\epsilon|a)$

3) 不具有任何两个连续a的所有符号串

$\epsilon|a|(b|ab)^+(\epsilon|a)$

若两个正规式U和V所表示的正规集相同, 则称为二者等价, 记为 $U=V$ .

例如,  $b(ab)^*=(ba)^*b$

$(a|b)^*=(a^*b^*)^*$

## 正规式的命名

■  $(0|1|2|\dots|9)(0|1|2|\dots|9)^*$



■  $<\text{digit}><\text{digit}>^*$

■ 其中,  $<\text{digit}> ::= 0|1|2|\dots|9$

■ 称,  $<\text{digit}> ::= 0|1|2|\dots|9$

■ 为名 $<\text{digit}>$ 的正规定义,

■ **regular definition of the name  $<\text{digit}>$**

11

## 正规定义

字母表 $\Sigma$ 上的正规定义形如:

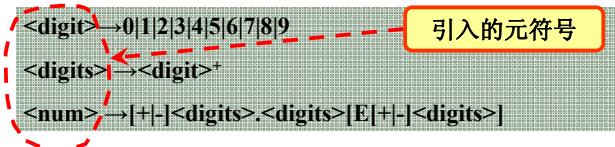
$d_1 \rightarrow r_1$

$d_2 \rightarrow r_2$

...

$d_n \rightarrow r_n$

各个 $d_i$ 的名字不同, 每个 $r_i$ 是字母表 $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$ 上的正规式, 其中 $i=1, 2, \dots, n$ .



12

### 正规式的代数性质

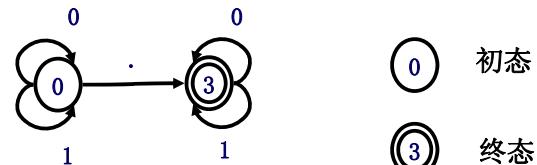
公理	描述
$U V = V U$	是可交换的
$U (V W) = (U V) W$	是可结合的
$(UV)W = U(VW)$	• 是可结合的
$U(V W) = UV UW$ $(U V)W = UW VW$	• 对 可分配
$\epsilon U = U$ $U \epsilon = U$	$\epsilon$ 是•的衡等元素
$U^{**} = U^*$	*是幂等的

13

- ⊕ 大写字母为正规式
- ⊕ | 正规式的或运算
- ⊕ • 正规式的连接运算
- ⊕ \* 正规式的闭包运算

### 3.1.2 Finite-State Automata

- Alphabet  $\Sigma$
- Set of states with initial and accept states
- Transitions between states, labeled with letters

 $(0|1)^*. (0|1)^*$ 

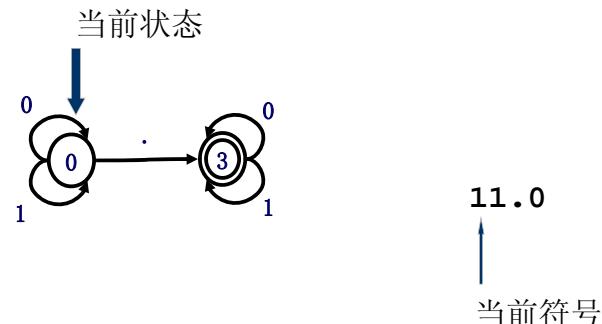
14

### Automaton Accepting String

- Conceptually, run string through automaton
- Have current state and current letter in string
- Start with start state and first letter in string
- At each step, match current letter against a transition whose label is same as letter
- Continue until reach end of string or match fails
- If end in accept state, automaton accepts string
- Language of automaton is set of strings it accepts

15

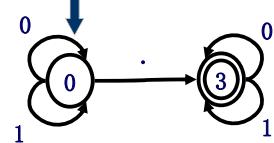
### 例



16

例

当前状态



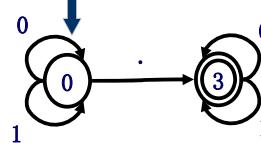
11.0

当前符号

17

例

当前状态



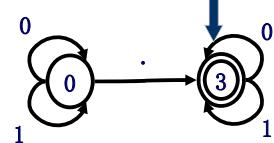
11.0

当前符号

18

例

当前状态



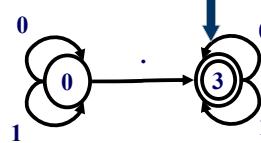
11.0

当前符号

19

例

当前状态



符号串被接受

11.0

当前符号

20

## 确定有限自动机(DFA)

一个确定有限自动机M是一个五元式

$M = (S, \Sigma, \delta, s_0, F)$ , 其中,

⊕ S是一个有限的**状态**集合;

⊕  $\Sigma$ 是有穷字母表, 它的每个元素称为一个**输入字符**;

⊕  $\delta$ 是一个从  $S \times \Sigma$  至  $S$  的单值部分映射。

$\delta(s, a) = s'$  意味着: 给定状态s和输入字符a返回状态s';

⊕  $s_0 \in S$ 是唯一的**初态**;

⊕  $F \subseteq S$ 是一个**终态集**(可空).

$$M = (\{0, 1, 2, 3\}, \{a, b\}, \delta, 0, \{3\})$$

21

一个DFA也可用一张(确定的)状态转换图来表示。假定DFA M含有m个状态和n个输入字符, 那么, 这个状态转换图含有m个状态结点, 每个结点顶多有n条箭弧射出和别的结点相连接, 整张图含有一个初态结点和若干个(可以为0)终态结点。

$$M = (\{0, 1, 2, 3\}, \{a, b\}, \delta, 0, \{3\})$$

状态	a	b
0	1	2
1	3	2
2	1	3
3	3	3

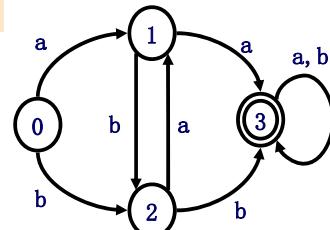


图 3.5 状态转换图

23

## 状态矩阵

一个DFA可用一个矩阵表示, 该矩阵的行表示状态, 列表示输入字符, 矩阵元素表示  $\delta(s, a)$  的值(转换的状态). 这个矩阵称为状态转换矩阵.

例: 有 DFA

$$M = (\{0, 1, 2, 3\}, \{a, b\}, \delta, 0, \{3\})$$

其中  $\delta$  为:

$$\delta(0, a) = 1 \quad \delta(0, b) = 2$$

$$\delta(1, a) = 3 \quad \delta(1, b) = 2$$

$$\delta(2, a) = 1 \quad \delta(2, b) = 3$$

$$\delta(3, a) = 3 \quad \delta(3, b) = 3$$

相应状态转换矩阵如下表: 其中

		a	b
0	1	2	
1	3	2	
2	1	3	
3	3	3	
		输入 符号	
		状态	
			矩阵元素

22

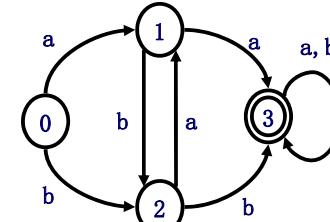
给定 DFA  $M = (S, \Sigma, \delta, s_0, F)$ , 对于  $\Sigma^*$  中的任意符号串  $\alpha$ , 若存在一条从初态结点到终态结点的路径, 且这条路径上的所有弧的标记依次连接成一个符号串, 同时这个符号串等于  $\alpha$  则称  $\alpha$  可以被 DFA M 所识别(或接受).

给定  $M = (\{0, 1, 2, 3\}, \{a, b\}, \delta, 0, \{3\})$  有:

$$S = \{0, 1, 2, 3\}$$

$$\Sigma = \{a, b\}$$

$$\Sigma^* = \{\epsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}$$



例如: 对于  $\Sigma^*$  上的符号串:

abaa 识别

所有为 M 所识别的  $\Sigma^*$  中的符号串组成的集合称为 M 的语言, 记为  $L(M)$ .

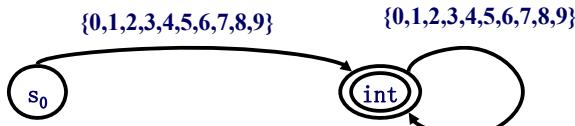
baba 不识别

24

## 正规式和DFA

如果一个DFA M的输入字母表为 $\Sigma$ , M称为 $\Sigma$ 上的一个DFA, 则 $\Sigma$ 上的一个字集 $V \subseteq \Sigma^*$ 是正规集, 当且仅当存在 $\Sigma$ 上的DFA M使得 $V=L(M)$ .

$(0|1|2|3|4|5|6|7|8|9) (0|1|2|3|4|5|6|7|8|9)^*$



$M = (\{s_0, \text{int}\}, \{0,1,2,3,4,5,6,7,8,9\}, \delta, s_0, \{\text{int}\})$

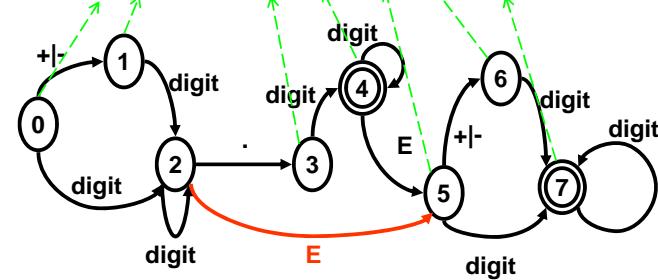
$\delta : \begin{pmatrix} \text{int} & \text{int} \\ \text{int} & \text{int} \end{pmatrix}$

25

$<\text{digit}> \rightarrow 0|1|2|3|4|5|6|7|8|9$

$<\text{digits}> \rightarrow <\text{digit}>^+$

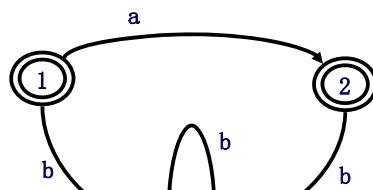
$<\text{num}> \rightarrow [+|-]<\text{digits}>.[<\text{digits}>|E|+|-]<\text{digits}>$



课堂练习: 写出该状态转换图所对应的DFA  
123E45 123.0 -123.0 +123.45 -123.45E-67

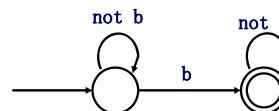
26

接受正规式  $ab^+ | ab^* | b^*$  的正规集的DFA如下所示:

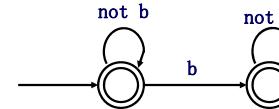


27

例: 只含一个b的串被如下所示DFA接受:



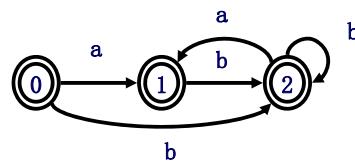
例: 包含最多一个b的串被如下所示的DFA接受:



28

## 例

- 不具有任何两个连续a的所有符号串



29

## 3.1.3 非确定有限自动机(NFA)

DFA  $M = (S, \Sigma, \delta, S_0, F)$  的确定性表现在  $\delta$  映射是一个单值函数. 即对于任何状态  $s \in S$ , 和输入符号  $a \in \Sigma$ ,  $\delta(s, a)$  唯一地确定了下一个状态.

当  $\delta(s, a)$  的值不唯一时...

30

## 3.1.3 非确定有限自动机(NFA)

一个非确定有限自动机  $M$  是一个五元式

$M = (S, \Sigma, \delta, S_0, F)$ , 其中,

- $S$  是一个有限的 **状态集合**.
- $\Sigma$  是有穷字母表, 它的每个元素称为一个 **输入字符**;
- $\delta$  是一个从  $S \times \Sigma^*$  至  $2^S$  的映射。  $\delta(s, \alpha) = T$  意味着: 给定状态  $s$  和输入符号串  $\alpha$  返回状态集合  $T$ , 其中  $T \subseteq S$ ;
- $S_0 \subseteq S$  是非空 **初态集**.
- $F \subseteq S$  是一个 **终态集(可空)**.

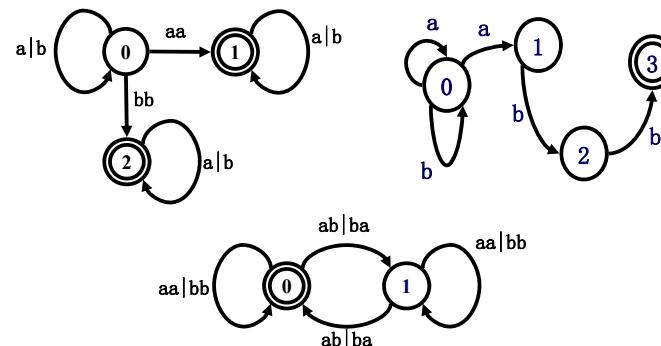
每个箭弧上的输入  
是一个符号串(含 $\epsilon$ )

同DFA

状态的迁移不再对于一个输入符号串是唯一的;  
开始状态可以多个

31

NFA 接受输入串  $\alpha$ , 当且仅当转换图中存在从开始状态到终态的一条路径, 使得该路径中每条边的输入符号串依次连接成为  $\alpha$



32

一个NFA也可用一张(确定的)状态转换图来表示。  
假定NFA M含有m个状态, 每个结点可射出若干个箭弧与别的结点相连, 每个弧用 $\Sigma^*$ 中的一个字作为输入字.  
整张图至少含有一个初态结点以及0到多个终态结点.

$$M = (\{0, 1, 2, 3\}, \{a, b\}, \delta, \{0\}, \{3\})$$

状态	a	b
0	{0,1}	{0}
1	$\varphi$	{2}
2	$\varphi$	{3}
3	$\varphi$	$\varphi$

$$\begin{pmatrix} \{0,1\} & \{0\} \\ \{\} & \{2\} \\ \{\} & \{3\} \\ \{\} & \{\} \end{pmatrix}$$

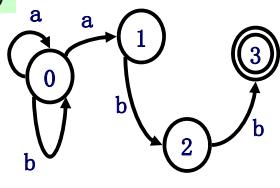


图 状态转换图  
33

对应的正规式为  $(a|b)^*abb$

NFA接受输入串 $\alpha$ , 当且仅当转换图中存在从开始状态到终态的一条路径, 使得该路径中每条边的输入字依次连接成为 $\alpha$ .

对应的正则式为  $aa^*|bb^*$

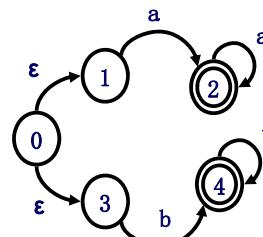
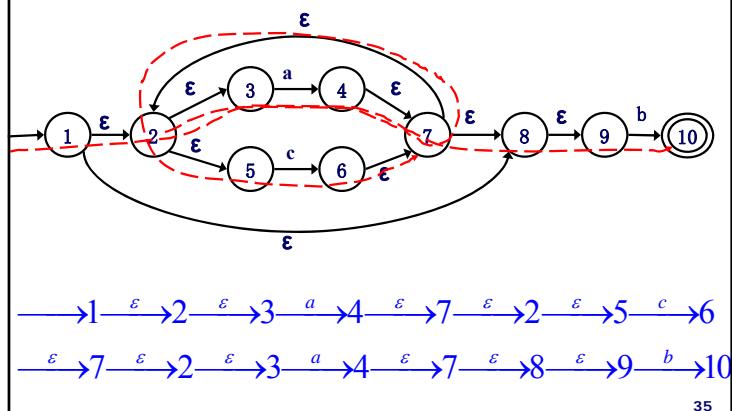


图 状态转换图  
34

34

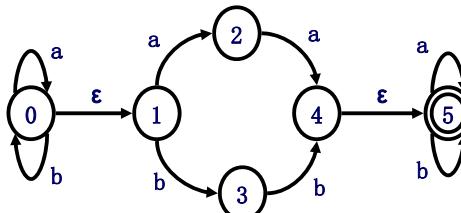
练习: 考虑以下NFA通过怎样的转换接受串acab:



$\xrightarrow{\epsilon} 1 \xrightarrow{\epsilon} 2 \xrightarrow{\epsilon} 3 \xrightarrow{a} 4 \xrightarrow{\epsilon} 7 \xrightarrow{\epsilon} 2 \xrightarrow{\epsilon} 5 \xrightarrow{c} 6$   
 $\xrightarrow{\epsilon} 7 \xrightarrow{\epsilon} 2 \xrightarrow{\epsilon} 3 \xrightarrow{a} 4 \xrightarrow{\epsilon} 7 \xrightarrow{\epsilon} 8 \xrightarrow{\epsilon} 9 \xrightarrow{b} 10$

35

NFA接受输入串 $\alpha$ , 当且仅当转换图中存在从开始状态到终态的一条路径, 使得该路径中每条边的输入字依次连接成为 $\alpha$ .



可接受的输入串为任意包含连续两个a或连续两个b的字符串  
 $\Sigma = \{a, b\}$

36

定理: 对于每个NFA M存在一个DFA M"使得 $L(M)=L(M'')$

证明: 假定NFA  $M = (S, \Sigma, \delta, S_0, F)$ , 对M的状态转换图进行改变, 构造等价的DFA  $M''$

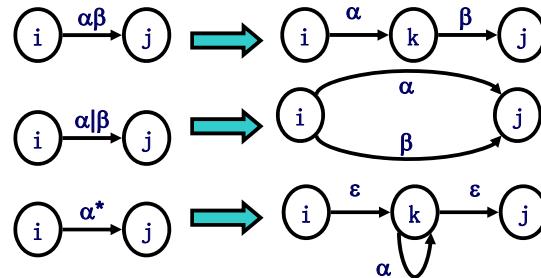
第一步: 由NFA M构造等价的NFA M'

第二步: 由NFA M'构造等价的DFA M''

37

①引进新的初态结点x和终态结点y, 从x到 $S_0$ 中每一状态结点连一条输入符号为 $\epsilon$ 的箭弧; 从F中每一状态结点连一条输入符号为 $\epsilon$ 的箭弧到y;

②重复进行如下图所示转换直到每条箭弧上的输入符号串的长度不大于1为止.



将最终得到的NFA记为M', 且有 $L(M)=L(M')$ .

38

### $\epsilon$ -CLOSURE(I) 的定义

DFA是NFA的特例, 可以采用子集法将NFA确定化.

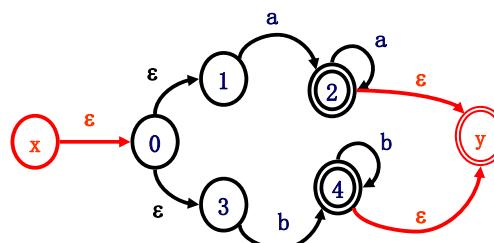
假定I是NFA M的状态集的一个子集, 我们定义 $\epsilon$ -CLOSURE(I)为:

- ⊕ 若 $s \in I$ , 则 $s \in \epsilon$ -CLOSURE(I);
- ⊕ 若 $s \in I$ , 那么从s出发经过任意条 $\epsilon$ 弧而能到达的任何状态 $s'$ 都属于 $\epsilon$ -CLOSURE(I).

状态集 $\epsilon$ -CLOSURE(I)称为I的 $\epsilon$ -闭包.

39

### 求 $\epsilon$ -CLOSURE(I)的例子



$$\epsilon\text{-CLOSURE}(\{x\}) = \{x, 0, 1, 3\}$$

$$\epsilon\text{-CLOSURE}(\{2\}) = \{2, y\}$$

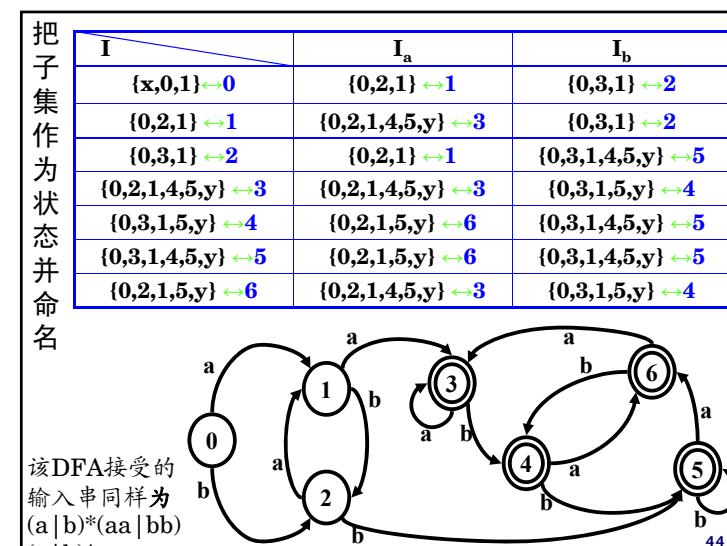
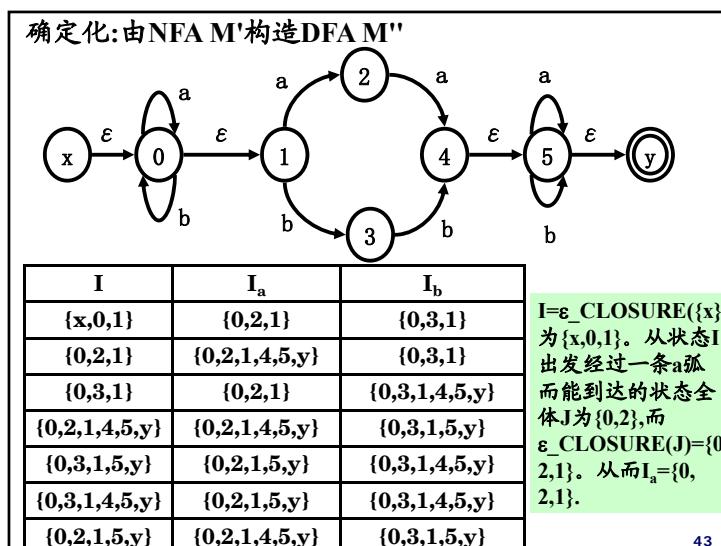
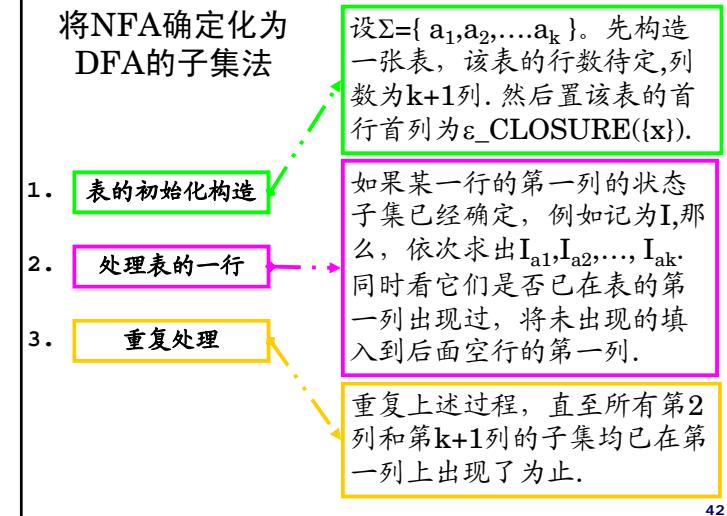
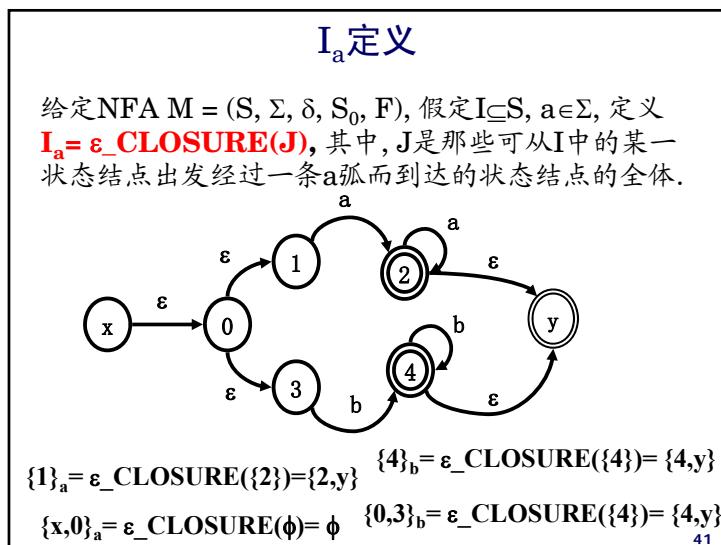
$$\epsilon\text{-CLOSURE}(\{x, 0, 2\}) = \{x, 0, 1, 3, 2, y\}$$

$$\epsilon\text{-CLOSURE}(\{1, 3, 4\}) = \{1, 3, 4, y\}$$

$$\epsilon\text{-CLOSURE}(\{0\}) = \{0, 1, 3\}$$

$$\epsilon\text{-CLOSURE}(\{4\}) = \{4, y\}$$

40

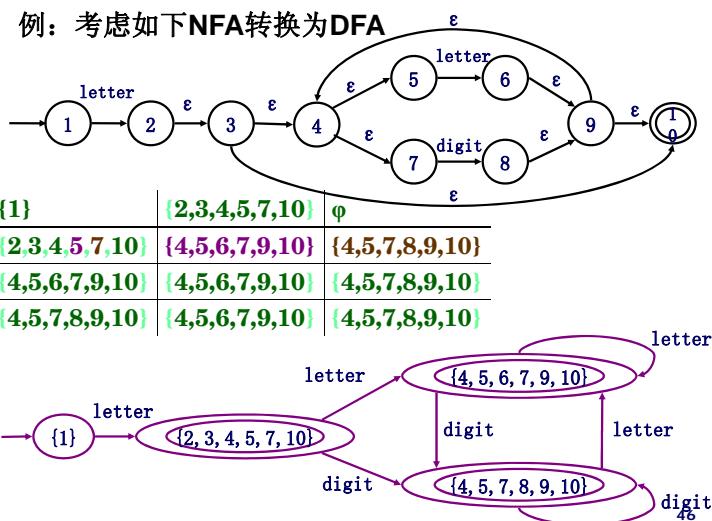


到此我们将这个表作为状态转换矩阵就得到

$$\text{DFA } M'' = (S', \Sigma', \delta'', s_0, F')$$

其中  $\delta''$  为通过子集算法构造的状态转换矩阵所表示的映射, 而  $s_0$  为  $\epsilon\text{-CLOSURE}(\{x\})$  命名后的状态名

45



46

### 3.1.4 3型文法与有限自动机的等价性

#### Chomsky文法

$$G = (V_t, V_N, S, \mathcal{P})$$

2型文法  $A \rightarrow \beta, A \in V_N, \beta \in (V_N \cup V_T)^*$

3型文法  
正规文法  $A \rightarrow \alpha B, \text{ 或, } A \rightarrow \alpha, \text{ 且 } A, B \in V_N, \alpha \in (V_T)^*$

47

#### 3型文法有两种表示形式

$$G = (V_t, V_N, S, \mathcal{P})$$

左线性文法:

$$A \rightarrow B\alpha, \text{ 或, } A \rightarrow \alpha, \text{ 且 } A, B \in V_N, \alpha \in V_t^*$$

右线性文法:

$$A \rightarrow \alpha B, \text{ 或, } A \rightarrow \alpha, \text{ 且 } A, B \in V_N, \alpha \in V_t^*$$

48

## 正规文法与有限自动机的等价性

对于3型文法 (**正规文法**)  $G$  和有限自动机  $M$ , 如果  $L(G)=L(M)$ , 则称  $G$  和  $M$  是等价的.

关于3型文法和有限自动机的等价性, 有以下结论:

- ① 对于每一个3型文法  $G$ , 都存在一个有限自动机  $M$ , 使得  $L(M)=L(G)$ .
- ② 对于每一个有限自动机  $M$ , 都存在一个3型文法  $G$ , 使得  $L(M)=L(G)$ .

49

**定理:** 对于每一个右线性文法  $G$ , 都存在一个有限自动机  $M$ , 使得  $L(M)=L(G)$ .

$G=(V_T, V_N, S, \mathcal{P})$  做如下构造得到的  $M$  为一个NFA

令  $M=(V_N \cup \{f\}, V_T, \delta, S, \{f\})$ , 其中  $f \notin V_N$ , 且  $\delta$  如下:

- 1) 若对于某个  $A \in V_N$  及  $a \in V_T \cup \{\epsilon\}$  有  $(A, a) \in \mathcal{P}$  则令  $\delta(A, a) = f$ ;
- 2) 对任意  $A \in V_N$  及  $a \in V_T \cup \{\epsilon\}$ , 设  $\mathcal{P}$  中左部为  $A$  右部第一符号为  $a$  的产生式形如  $A \rightarrow aA_1 | aA_2 | \dots | aA_k$  则令  $\delta(A, a) = \{A_1, \dots, A_k\}$ .

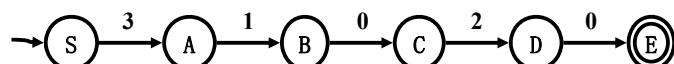
50

## 举例: 自动机与右线性文法

■ 符号串 31020

■ 部分规则

$$\begin{aligned} S &\rightarrow 3A \\ A &\rightarrow 1B \\ B &\rightarrow 0C \\ C &\rightarrow 2D \\ D &\rightarrow 0 \end{aligned}$$



■  $S \Rightarrow 3A \Rightarrow 31B \Rightarrow 310C \Rightarrow 3102D \Rightarrow 31020$

51

**定理:** 对于每一个有限自动机  $M$ , 都存在一个右线性文法  $G$ , 使得  $L(M)=L(G)$ .

DFA  $M = (S, \Sigma, \delta, s_0, F)$  做如下构造得到右线性文法  $G$

- 1) 若  $s_0 \notin F$ , 令  $G = (\Sigma, S, s_0, \mathcal{P})$ , 其中  $\mathcal{P}$  构造如下:

对于任何  $a \in \Sigma$  及  $A, B \in S$ , 若有  $\delta(A, a) = B$ , 则:

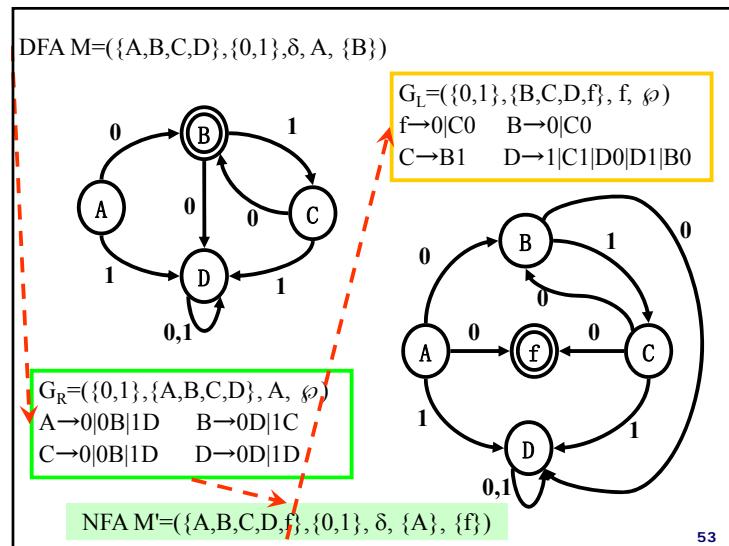
- i) 当  $B \notin F$  时, 令  $A \rightarrow aB$
- ii) 当  $B \in F$  时, 令  $A \rightarrow a | aB$

2)  $s_0 \in F$ , 因为  $\delta(s_0, \epsilon) = s_0$ , 所以  $\epsilon \in L(M)$ , 从而  $L(G) = L(M) - \{\epsilon\}$ .

因此令  $G' = (\Sigma, S \cup \{s'_0\}, s'_0, \mathcal{P} \cup \{(s'_0, s_0), (s'_0, \epsilon)\})$

其中  $s'_0 \notin S$ . 最后令  $G = G'$

52

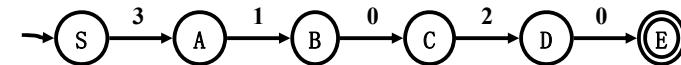


### 举例：自动机与左线性文法

■ 符号串 31020

■ 部分规则

$A \rightarrow 3$   
 $B \rightarrow A1$   
 $C \rightarrow B0$   
 $D \rightarrow C2$   
 $E \rightarrow D0$



■  $E \Rightarrow D0 \Rightarrow C20 \Rightarrow B020 \Rightarrow A1020 \Rightarrow 31020$

54

### 3.1.5 正规式与有限自动机的等价性

关于正规式和有限自动机的等价性，有以下定理：

定理：对于字母表 $\Sigma$ 上的任何有限自动机 $M$ ，都存在一个 $\Sigma$ 上的正规式 $r$ ，使得 $L(r) = L(M)$ . 反之亦然。

构造性证明：

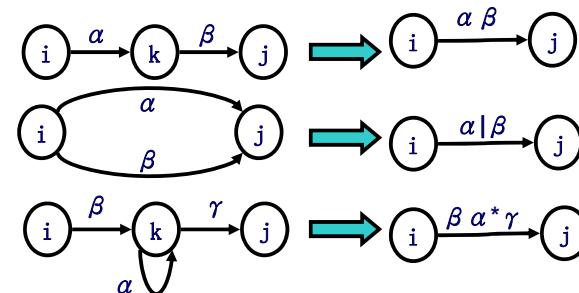
1. 对 $\Sigma$ 上的NFA  $M$ 构造 $\Sigma$ 上的正规式  $r$  使得 $L(M) = L(r)$ .

2. 对 $\Sigma$ 上的正规式  $r$ , 构造 $\Sigma$ 上的NFA  $M$ , 使得 $L(r) = L(M)$ .

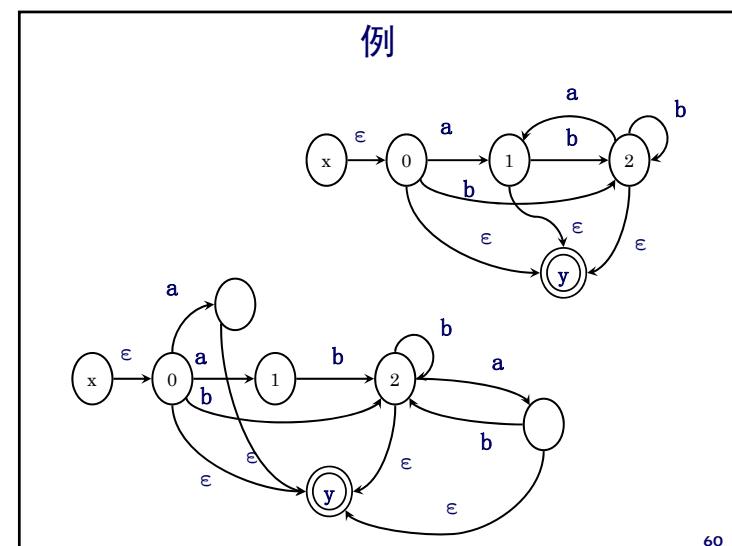
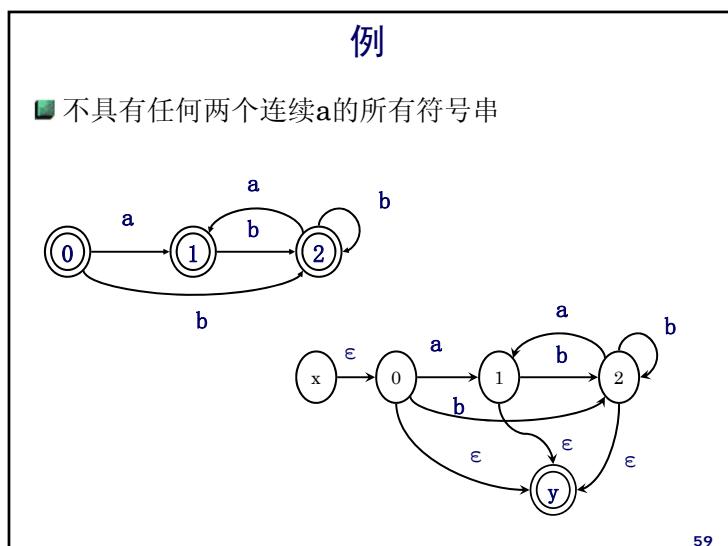
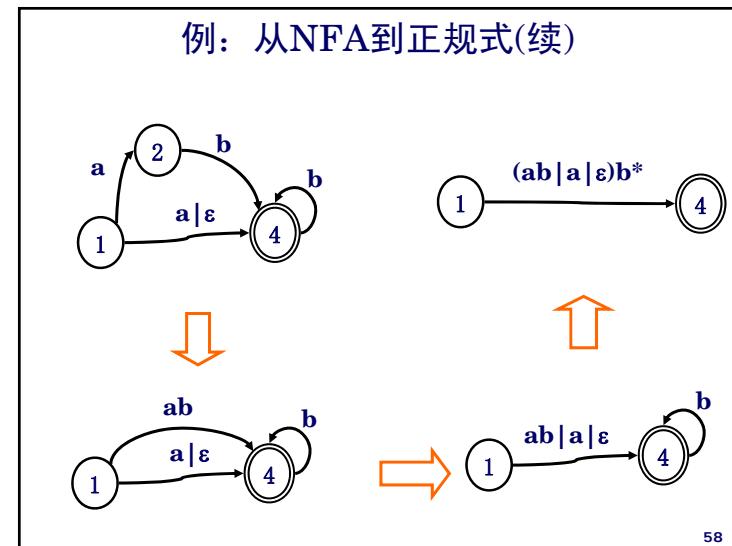
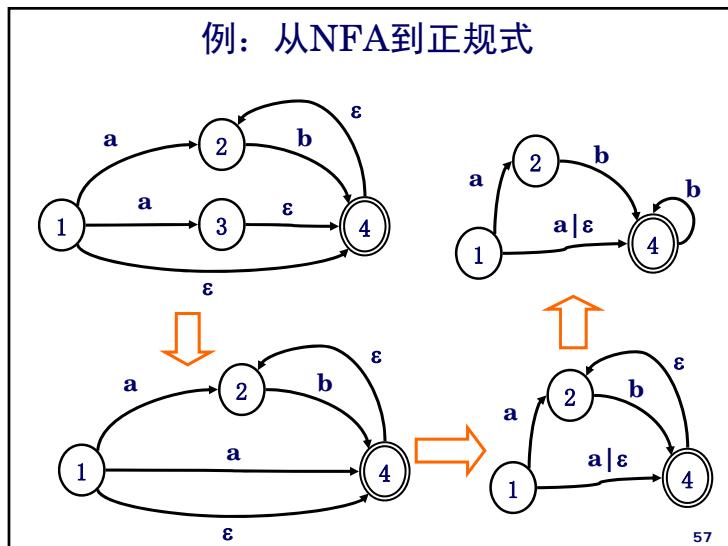
55

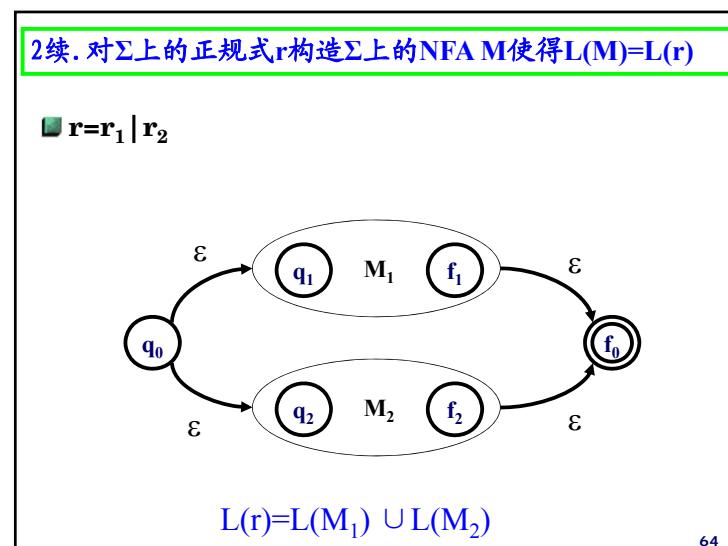
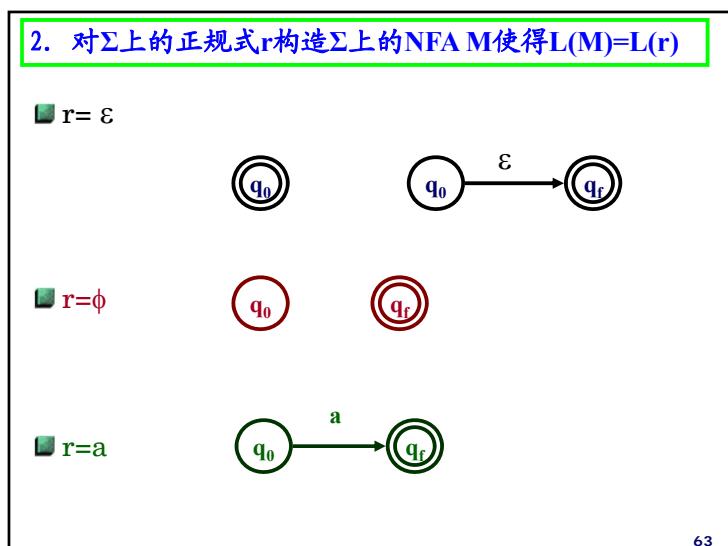
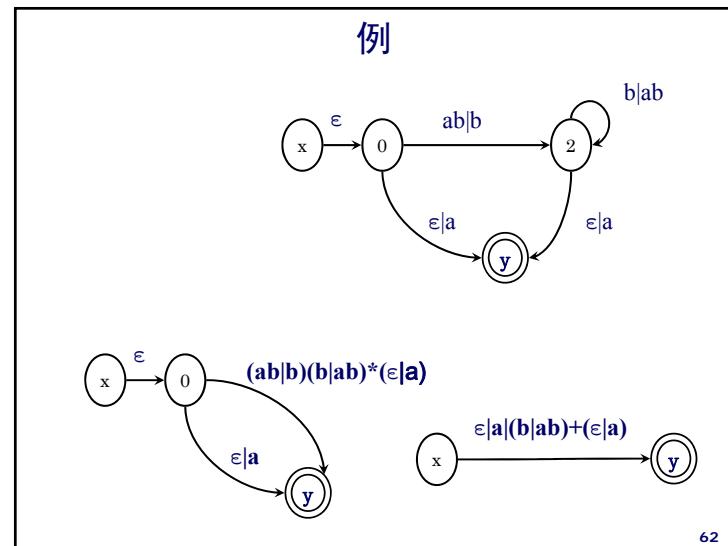
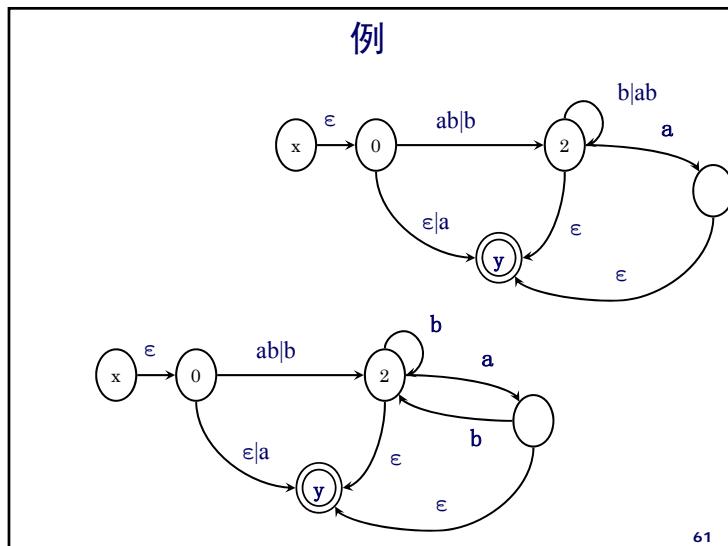
令每条弧以正规式作标记；在图中加进X和Y两个结点，从X用 $\epsilon$ 弧连接到M的所有初态结点，从M的所有终态结点用 $\epsilon$ 弧分别连接到Y；将X作为初态，Y作为终态构成新的一个NFA  $M'$ 且 $L(M) = L(M')$ .

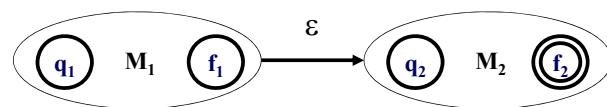
采用图示替换规则消去 $M'$ 中所有结点，只剩下X和Y，则X到Y的弧的标记就是R



56

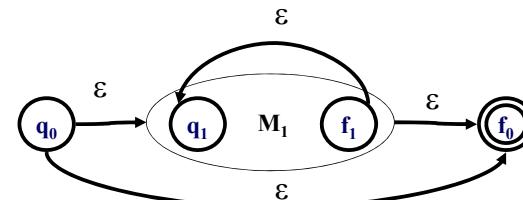




2续. 对 $\Sigma$ 上的正规式 $r$ 构造 $\Sigma$ 上的NFA  $M$ 使得 $L(M)=L(r)$   $r=r_1r_2$ 

$$L(r)=L(M_1) L(M_2)$$

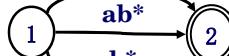
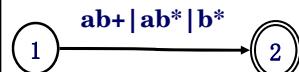
65

2续. 对 $\Sigma$ 上的正规式 $r$ 构造 $\Sigma$ 上的NFA  $M$ 使得 $L(M)=L(r)$   $r=r_1^*$ 

$$L(r)=L(M_1)^*$$

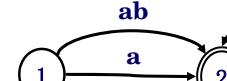
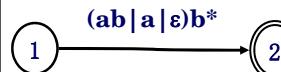
66

## 例：从正规式到 NFA

  $ab+ | ab^* | b^*$ 

67

## 例：从正规式到 NFA (续)

  $ab+ | ab^* | b^*$ 

68

## 结论

- 正规文法、正规式、确定有限自动机、非确定有限自动机在接受语言能力上等价。

69

## 3.1.6 确定有限自动机的化简

一个DFA M的化简是指：

寻找一个状态数比M少的DFA M'，使得 $L(M)=L(M')$ 。

一个DFA M的状态最少化过程旨在将M的状态集分割成一些不相交的子集，使得任何不同的两个子集中的状态都是可区别的，而同一子集中的任何两个状态都是等价的。最后，在每个子集中选出一个代表，同时消去其它等价状态。

M的两个状态s和t是等价的：当且仅当分别从二者出发至终态均能够读出同一字。

M的两个状态s和t是可区别的：当且仅当二者不等价。

70

### 状态集S的划分步骤：

- (1) 基本划分{终结状态集, 其他状态集}是一个划分 $\Pi$ ；
- (2) 对于当前划分 $\Pi=\{I^{(1)}, \dots, I^{(m)}\}$ ,  $m \geq 1$ ; 进行如下处理：  
对于每个 $a \in \Sigma$ , 若 $I_a^{(i)}$ 分别与 $\Pi$ 的 $n$ 个不同块相交不为空 (其中 $1 < n \leq m$ )，则将 $I^{(i)}$ 划分成 $n$ 个不相交的子集，使得每个子集 $J$ 的 $J_a$ 包含于 $\Pi$ 的某个块。然后将每个 $J$ 加入 $\Pi$ 并删除 $I^{(i)}$
- (3) 重复上述过程(2)直至 $\Pi$ 不再变化为止。

$$\Pi = \{X_1, \dots, X_m\}, \quad X_i \subseteq X, i=1 \dots m, m \geq 1$$

$$X_i \cap X_j = \emptyset, \bigcup_{i=1}^m X_i = X$$

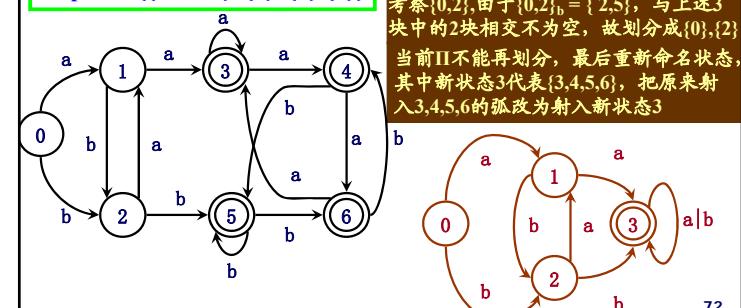
71

### 例：DFA的化简：

Step1:  $\Pi = \{\{3,4,5,6\}, \{0,1,2\}\}$

Step2:  $\Pi = \{\{3,4,5,6\}, \{1\}, \{0,2\}\}$

Step3:  $\Pi = \{\{3,4,5,6\}, \{1\}, \{0\}, \{2\}\}$



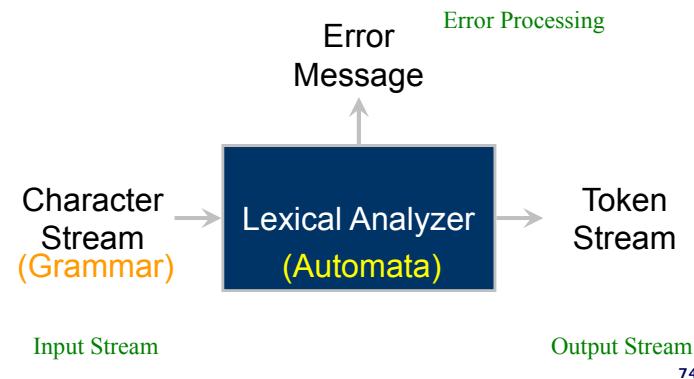
72

### 3.2 词法分析器的设计与实现

- Introduction
- Design Issues
- Examples
- Generators

73

#### 3.2.1 Introduction to Lexical Analyzer



74

#### Character Stream

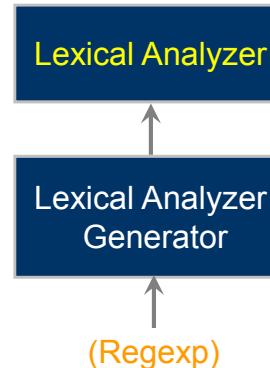
- Source Code
  - ⊕ Programming Languages
  - ⊕ Script Languages (HTML, XML)
  - ⊕ Others
- Grammar
  - ⊕ Alphabet
  - ⊕ Production rules

#### Token Stream

- Category
- Value
- Symbol Table
- Other Tables

75

#### Lexical Analyzer Generator



76

### 3.2.2 Design Issues on Lexical Analyzer

- Token的表示
  - 全体一种
  - 一字一种
- 词法分析器与语法分析器的关系
  - 单独一遍
  - 作为语法分析器的子程序
- 输入缓冲区的控制
  - 预处理
  - 超前搜索
- DFA的表示
  - Matrix
  - Case语句

77

### (1) 单词符号的分类

- **关键字**: 由程序语言定义的具有固定意义的标识符。也可称为保留字或基本字。例如: Pascal中的begin, end, if等。它是确定的。
- **标识符**: 用来表示各种名字, 如变量名、数组名、过程名等。它是不限的。
- **常数**: 常数的类型一般有整型、实型、布尔型、文字型等。它是不限的。
- **运算符**: 如+、-、\*、/等。它是确定的。
- **界符**: 如逗号、分号、括号、/\*, \*/等。它是确定的。

78

### (1) Token的表示

- 用二元组表示: <种别, 属性值>
- **种别**: 采用整数表示、常量表示或者特殊标识符表示
- **属性值**: 对于常数采用内部表示; 标识符(变量、常量、数组名、函数名等)采用符号表入口, 等等。

**种别的确定:** 主要取决于实现上的方便, 例如: 标识符一般同归为一种。常数宜按类型(整、实、布尔)分。关键字可以将其全体视为一种, 也可一字一种。运算符可采用一符一种, 但也可把具有一定共性的视为一种。界符则一般采用一符一种。若是一符一种分种, 单词自身值就不需要了。

79

### 例

- 例: FORTRAN编译程序的词法分析器在扫描输入串 IF (5·EQ·M) GOTO 100 后, 它输出的**单词符号流**是:

逻辑IF	(34, _)
左括号	(2, _)
整常数	(20, ‘5’的二进制表示)
等号	(6, _)
标识符	(26, ‘M’)
右括号	(16, _)
GOTO	(30, _)
标号	(19, ‘100’的二进制表示)

80

### 例

- 例：C++代码段：**while ( i >= j ) i --;** 经词法分析器处理以后，它将被转换为如下的单词符号流：

```
( while , _ )
( ( , _ )
( id , 指向i的符号表指针 )
( >= , _ )
( id , 指向j的符号表指针 )
( ) , _ )
( id , 指向i的符号表指针 )
( -- , _ )
( ; , _ )
```

81

### (2)词法分析器与语法分析器的关系

- 作为不同的遍
- 作为子程序调用

82

### (3)超前搜索问题

- 例：关键字识别中的超前搜索：

例如：在标准FORTRAN中

- 1、DO99K = 1,10      \_\_\_\_\_ 其中的DO、
- 2、IF(5.EQ.M)I = 10    \_\_\_\_\_ IF为关键字
- 3、DO99K = 1.10       \_\_\_\_\_ 其中的DO、
- 4、IF(5) = 55          \_\_\_\_\_ IF为标识符的一部分

83

#### ■ 标识符的识别

- ◆ 多数语言的标识符是字母开头的“字母/数字”串，而且在程序中标识符的出现后都跟着算符或界符。因此，不难识别。

#### ■ 常数的识别

- ◆ 对于某些语言的常数的识别也需要使用超前搜索。

#### ■ 算符和界符的识别

- ◆ 对于诸如C++语言中的“++”、“--”，这种复合成的算符，需要超前搜索。

84

#### (4) 状态转换图的实现

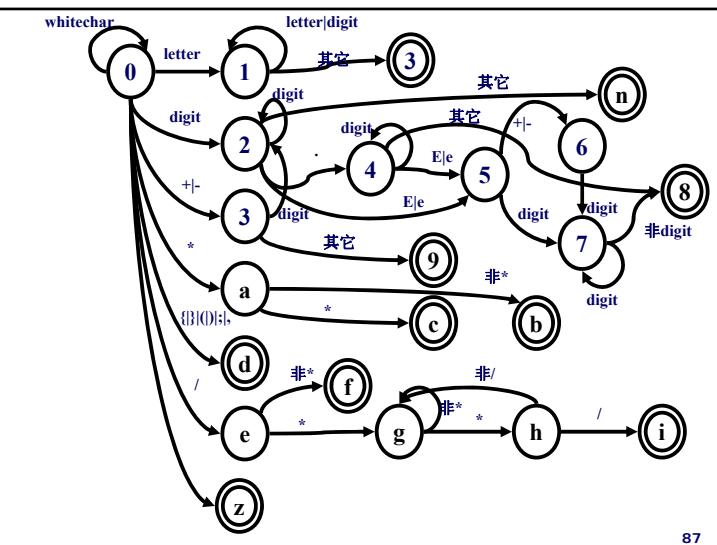
1. ch 字符变量, 存放最新读进的源程序字符。
2. TOKEN 字符数组, 存放构成单词的字符串。
3. GETCHAR 过程, 将下一输入字符读入ch, 搜索指示器前移一个字符。
4. GETBC 过程, 检查ch中的字符是否为空白。若是, 则调用GETCHAR 直至CHAR中进入一个非空白字符。
5. CONCAT 过程, 把ch中的字符连接到TOKEN之后。
6. IsLetter 布尔函数过程, 它们分别判断ch中的字符是数字或是字母, IsDigit 从而给出真假值TRUE、FALSE。
7. Reserve 整型函数过程, 用TOKEN中的字符串查保留字表, 若是一个保留字则给予编码, 否则回送0值(假定0不是保留字的编码)。
8. Retract 过程, 把搜索指示器回调一个字节, 把ch中的字符置为空白。
9. InsertID 插入符号表, 返回入口(指针)
10. InsertConst 插入常数表, 返回入口(指针)

85

#### 3.2.3 Examples

##### ■ C程序 (作为子程序)

86



87

#### 符号表

```
struct entry
{
    struct entry *next; /* linked list implementation */
    char *name; /* the user-defined name of the identifier */
    int type; /* its type */
    int blockno; /* scoping information */
    int addr; /* address assigned by code generator */
};
```

88

```
/* constants for our scanner */
#define LPAREN '('          /* single char tokens */
#define RPAREN ')'          /* single char tokens */
#define COMMA ','           /* single char tokens */
#define GREATER '>'         /* single char tokens */
#define LESS '<'            /* single char tokens */
#define ASSIGN '='           /* single char tokens */
...
#define WHILE 257           /* reserved words */
#define PRINT 258           /* reserved words */
#define RETURN 259          /* reserved words */
#define IF 260               /* reserved words */
...
#define VARIABLE 268        /* other tokens */
#define INTEGER 269          /* other tokens */
#define TEXT 270             /* other tokens */
#define EOF 271              /* end of file */
#define UNKNOWN 272          /* don't recognize */
```

```
/* some globals */

int ch;                      /* next char */
int intval;                  /* integer lexeme */
char textval[255];           /* text string lexeme */

void init()
{
    /* set up a lookup table for reserved words */
    create_reserved_table();
    insert_reserved(WHILE, "WHILE")
    insert_reserved(PRINT, "PRINT")
    insert_reserved(RETURN, "RETURN")
    ...
    /* get the first character */
    ch = getchar();
}
```

90

```
int scanner(){
switch(ch){
    case ' ': /* white space */
    case '\n':
    case '\t':
        while (isspace(ch = getchar())); /* get rid of the white */
        return scanner();
    case '/': /* check for a comment or DIVIDE operator */
        ch = getchar();
        if (ch != '/')
            return DIVIDE;
        ... /* code to skip over a comment */
    case '(':
    case ')':
    case ',':... /* and any other single char tokens */
        intval = ch;
        ch = getchar();
        return intval;
    ...
```

91

```
case 'A': /*assume reserved words are caps, variables lower case*/
case 'B':
case 'C': ... /* reserved words */
    textval[0] = ch;
    for (i = 0; isupper(ch = getchar()); textval[i++] = ch)
        ;
    textval[i++] = '\0';
    return lookup_reserved(textval);
case 'a':
case 'b':
case 'c':... /* variables */
    textval[0] = ch;
    for (i = 0; islower(ch = getchar()); textval[i++] = ch);
    textval[i++] = '\0';
    if (lookup_symtab(textval) == UNKNOWN)
        add_symtab(textval);
    return VARIABLE;
```

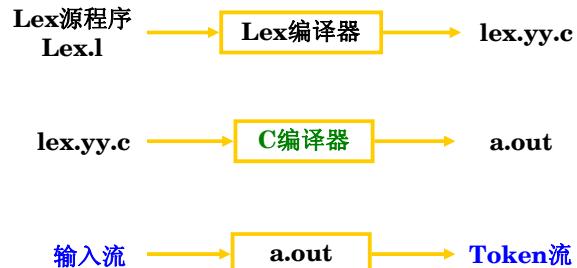
```

case '0':
    case '1':....          /* integer */
        intval = ch - '0'; /* convert to a number */
        while (isdigit(ch = getchar()))
            intval = intval * 10 + ch - '0';
        return INTEGER;
case EOF:
    return EOF;
default:
    intval = ch;
    ch = getchar();
    return UNKNOWN;
}

```

93

### 3.2.4 词法分析器的自动产生



94



《lex与yacc (第二版)》  
 作者: [John R. Levine](#), [Tony Mason](#), [Doug Brown](#) 著 杨作梅, 张旭东, 等 译  
 出版: 2003年1月  
 书号: 7-111-10721-7  
 页数: 392  
 定价: 45.00元

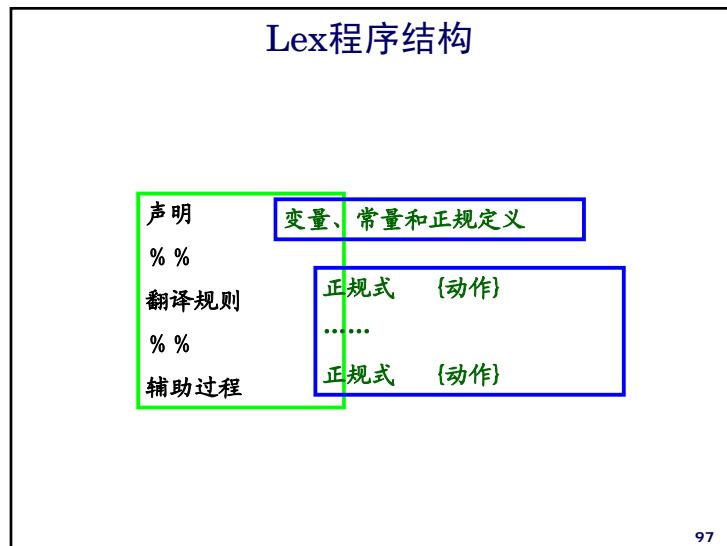
95

```

<程序> ::= <程序首部>; <分程序>。
<程序首部> ::= program <标识符>
<分程序> ::= <复合语句>
<复合语句> ::= begin <语句序列> end
<语句序列> ::= <语句> {; <语句>}
<语句> ::= <赋值语句> | <复合语句> | <条件语句>
<赋值语句> ::= <标识符> := <表达式>
<条件语句> ::= if <布尔表达式> then <语句> else <语句>
<表达式> ::= <项> { (+|-) <项> }
<项> ::= <因式> { (*|/) <因式> }
<因式> ::= <标识符> | <无正负号常量> | ' (' <表达式> ')'
<布尔表达式> ::= <表达式> <关系运算符> <表达式>
<关系运算符> ::= = | <| <= | > | >= | <>
<标识符> ::= <字母> {<字母> | <数字>}
<无正负号常量> ::= <数字> {<数字>} [ . <数字> {<数字>} ]
<字母> ::= a | b | c | d | e | f | g | ..... | u | v | w | x | y | z
<数字> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

```

输入以“#”为结束符



A-Z, 0-9, a-z	构成了部分模式的字符和数字。
.	匹配任意字符，除了 \n。
-	用来指定范围。例如：A-Z 指从 A 到 Z 之间的所有字符。
[]	一个字符集合。匹配括号内的任意字符。如果第一个字符是 ^ 那么它表示否定模式。例如: [abc] 匹配 a, b, 和 c 中的任何一个。
*	匹配 0 个或者多个上述的模式。
+	匹配 1 个或者多个上述模式。
?	匹配 0 个或 1 个上述模式。
\$	作为模式的最后一个字符匹配一行的结尾。
{}	指出一个模式可能出现的次数。例如: A{1,3} 表示 A 可能出现1次或3次。
\	用来转义元字符。同样用来覆盖字符在此表中定义的特殊意义，只取字符的本意。
^	否定。
	表达式间的逻辑或。
"<一些符号>"	字符的字面含义。元字符具有。
/	向前匹配。如果在匹配的模版中的 "/" 后跟有后续表达式，只匹配模版中 "/" 前面的部分。如：如果输入 A01，那么在模版 A0/1 中的 A0 是匹配的。
( )	将一系列常规表达式分组。

98

```

%{常量定义... %}
delim   [ \t\n]
ws       {delim}+
letter   [A-Za-z]
digit   [0-9]
id      {letter}{digit}*
number  {digit}+(\.{digit}+)?(E[+\-]?{digit}+)?

%%
{ws}  {}
if     {return(IF);}
then   {return(THEN);}
else   {return(ELSE);}
{id}   {yyval=install_id();return(ID);}
{number} {yyval=install_num();return(NUMBER);}
“<”   {yyval=LT;return(RELOP);}

%%
install_id() { 单词装入符号表并返回指针； }
...
    
```

99

### 上机实验I

- 见论坛上的课程上机题目，根据给定的C<sub>0</sub>语言的文法，采用Lex或者编写程序构造词法分析器。
- 上机时间第四周、第五周：三个班安排在同一时间。由班级负责同学到计算机系实验中心登记，定了以后通知辅导老师。
- 实验部分以当时上机情况和提交的报告进行评分，占本课程成绩10%。

100

## 实验题目

- 实验I: 构造词法分析器, 采用Lex生成, 或者编程实现。
- 实验II: 构造语法分析器, 采用Yacc生成, 或者编程实现。
- 实验III\*: 产生中间代码, 具体内容另给。

在论坛上公布上机文档:

- 实验I: 借助Lex和Yacc进行词法语法分析
- 实验一lex与yacc的使用. doc
- 实验II: 用C或Java编写词法分析器
- 实验二词法分析器. doc
- 实验III\*: 用C或Java编写语法分析器
- 实验三语法分析. doc

101

## C<sub>0</sub>语法

- ⊕ program → { var-declaration | fun-declaration }
- ⊕ var-declaration → int ID { , ID }
- ⊕ fun-declaration → ( int | void ) ID ( params ) compound-stmt
- ⊕ params → int ID { , int ID } | void | empty
- ⊕ compound-stmt → { { var-declaration } { statement } }
- ⊕ statement → expression-stmt | compound-stmt | if-stmt | while-stmt | return-stmt
- ⊕ expression-stmt → [ expression ] ;

102

## C<sub>0</sub>语法 (续)

- ⊕ if-stmt → if( expression ) statement [ else statement ]
- ⊕ while-stmt → while( expression ) statement
- ⊕ return-stmt → return [ expression ] ;
- ⊕ expression → ID = expression | simple-expression
- ⊕ simple-expression → additive-expression [ relop additive-expression ]

103

## C<sub>0</sub>语法 (续)

- ⊕ relop → < | <= | > | >= | == | !=
- ⊕ additive-expression → term [ ( + | - ) term ]
- ⊕ term → factor [ ( \* | / ) factor ]
- ⊕ factor → ( expression ) | ID | call | NUM
- ⊕ call → ID( args )
- ⊕ args → expression { , expression } | empty
- ⊕ ID → ... ;参见C语言标识符定义
- ⊕ NUM → ... ;参见C语言数的定义

104

### 报告提交

- 实验报告应该按照格式填写。
- 报告中写明程序设计思想及简要设计方案，代码简要描述，运行测试情况，体会。
  
- 每人一共提交2个报告，对应实验I和II（含III可选）
- 提交方式：作为附件email到指定邮箱，若得到回复确认表示提交成功。
- 指定邮箱为：[yinliang\\_zhao@163.com](mailto:yinliang_zhao@163.com)
- 邮件标题格式：姓名-学号-班级-实验报告n
- 附件命名格式：姓名-学号-班级-实验报告n.rar
- 附件(.rar)中含有：源程序、测试数据结运行结果、报告、其他辅助文件（若有）

105

### 本章练习

- 作业27日交：P63-65:2,6-9,12,14,15
- 上机：完成上机题I并写好报告

106

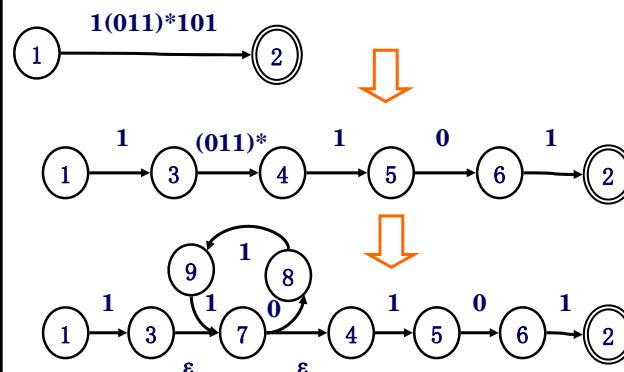
### P63 习题6 (2)

- A是任意正规式，证明 $(A^*)^* = A^*$
- 只须证明 $L((A^*)^*) = L(A^*)$ 
  - ⊕ 对任意 $\alpha \in L(A^*)$ ,  $\alpha = A^k$ , 其中k为非负整数;
  - ⊕ 那么 $\alpha$ 可写成 $(A^k)^1$ , 即与正规式 $(A^*)^*$ 匹配;
  - ⊕ 因此,  $\alpha \in L((A^*)^*)$ .
- 对任意 $\alpha \in L((A^*)^*)$ ,  $\alpha = (A^k)^j$ , 其中k和j为非负整数;
  - ⊕ 那么 $\alpha$ 可写成 $A^{kj}$ , 其中 $k*j$ 为非负整数;
  - ⊕ 所以 $\alpha$ 与正规式 $A^*$ 匹配, 从而,  $\alpha \in L(A^*)$ .

107

### P64 习题7 (1)

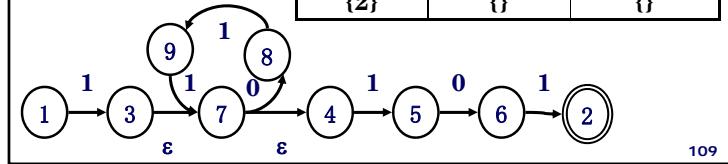
- 构造下列正规式相应的DFA,  $1(011)^*101$



108

P64 习题7 (1)续

I	I <sub>0</sub>	I <sub>1</sub>
{1}	{}	{3,7,4}
{3,7,4}	{8}	{5}
{8}	{}	{9}
{5}	{6}	{}
{9}	{}	{7,4}
{6}	{}	{2}
{7,4}	{8}	{5}
{2}	{}	{}



P64 习题7 (1)续

I	I <sub>0</sub>	I <sub>1</sub>
{1}	1	{3,7,4}
{3,7,4}	Q	{8}
{8}	8	{5}
{5}	5	{6}
{9}	9	{7,4}
{6}	6	{2}
{7,4}	R	{8}
{2}	2	{}

