



第3章

词法分析

本章主要内容

- 词法分析程序
- 单词的形式化描述工具
- 有穷自动机
- 有穷自动机和正规表达式
- 有穷自动机和正规文法
- 词法分析程序的自动构造

回顾：什么是词法分析程序？

习题P64页1(1), 4(a), 5, 8

3.1 词法分析程序的设计

3.1.1 词法分析程序和语法分析程序接口方式

- 可以**作为**单独的一遍
- 较常用的方式是由**语法分析程序**调用
- 基本任务都是**识别单词**

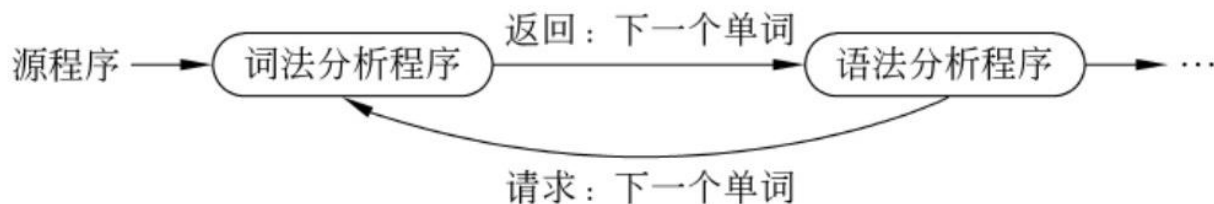


图 3.1 语法分析程序调用词法分析程序

3.1.2 词法分析程序的输出

- 词法分析程序的功能是读入源程序，输出单词符号。
- 词法分析程序所输出的单词符号常常用以下二元式表示：

（单词种别，单词自身的值）

- 单词的种别是语法分析需要的信息
- 单词的值则是编译其它阶段需要的信息

■ 程序语言单词的分类:

1. 标识符：用来表示各种名字
2. 字面常数：256, 3.14, true, 'abc'
3. 关键字(保留字或基本字): begin, end
4. 运算符：如, +、-、*、/ 等等
5. 界符：如逗号, 分号, 冒号等

单词种类用整数编码表示:

■ 例

if i=5 then x:=y;

保留字 if	(3, 'if')
标识符 i	(1, 指向 i 的符号表入口)
等号 =	(4, '=')
常数 5	(2, '5')
保留字 then	(3, 'then')
标识符 x	(1, 指向 x 的符号表入口)
赋值号 :=	(4, ':=')
标识符 y	(1, 指向 y 的符号表入口)
分号 ;	(5, ';')

PL/0编译程序中词法分析程序

■ 文法

<无符号整数> ::= <数字> {<数字>}
<标识符> ::= <字母> {<字母> | <数字>}
<字母> ::= *a* | *b* | ... | *X* | *Y* | *Z*
<数字> ::= 0 | 1 | 2 | ... | 8 | 9
<保留字> ::= const | var | procedur | begin | end | odd
 | if | then | call | while | do | read | write
<运算符> ::= + | - | * | / | = | # | < | <= | > | >= | :=
<界符> ::= (|) | , | ; | .

词法单位 (31个单词种别)	标识符 1 个	: ident
	无符号整数 1 个	: number
	保留字 13 个	: plus minus ...
	运算符 11 个	: beginsym endsym ...
	界符 5 个	: lparen rparen ...

PL/0二元组示例

- 输出表示（单词种别，单词自身的值）。

A=B+2

(ident,指向A的符号表的入口指针)

(becomes,)

(ident,指向B的符号表的入口指针)

(plus,)

(number, 2)

词法分析程序

- 主要任务:
 - 读源程序，产生单词符号
- 其他任务:
 - 滤掉空格，跳过注释、换行符
 - 追踪换行标志，复制出错源程序，
 - 宏展开， ...

* 词法分析, 获取一个符号

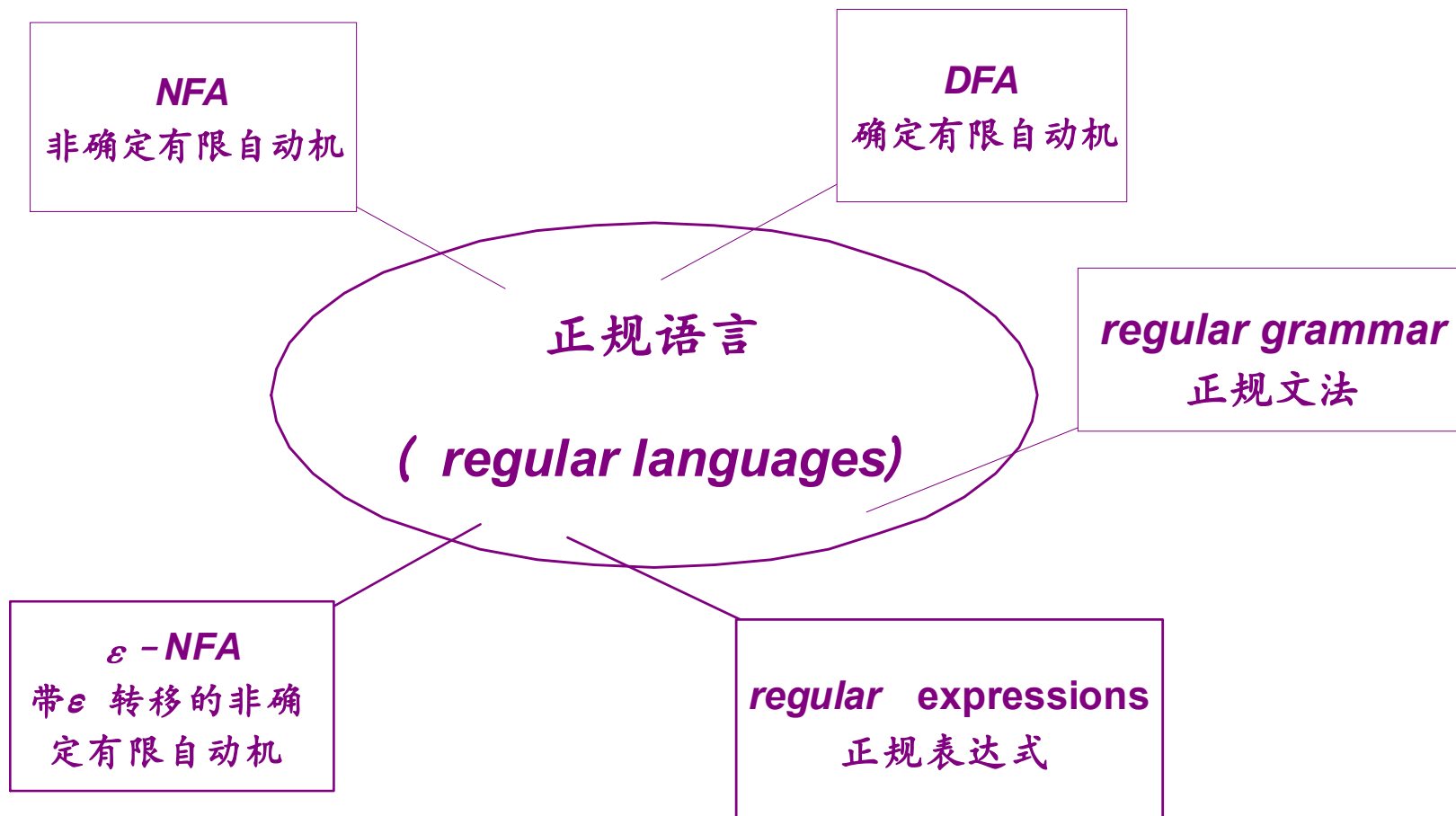
*/

```
int getsym()
{
    int i,j,k;

    /* the original version lacks "\r", thanks to foolevery */
    while (ch==' ' || ch==10 || ch==13 || ch==9) /* 忽略空格、换行、回车和TAB */
    {
        getchdo;
    }
    if (ch>='a' && ch<='z')
    {
        /* 名字或保留字以a..z开头 */
        k = 0;
        do {
            if(k<al)
            {
                a[k] = ch;
                k++;
            }
            getchdo;
        } while (ch>='a' && ch<='z' || ch>='0' && ch<='9');
        a[k] = 0;
        strcpy(id, a);
        i = 0;
        j = norw-1;
        do {
            /* 搜索当前符号是否为保留字 */
            k = (i+j)/2;
            if (strcmp(id,word[k]) <= 0)
            {
                j = k - 1;
            }
            if (strcmp(id,word[k]) >= 0)
            ,
        }
```

- 保留字
- 运算符
- 标识符
- 无符号整数
- 界符

3.3 单词的形式化描述工具



3.3.1 正规文法

对于文法中每一条产生式 $\alpha \rightarrow \beta$ 的形式都为 $A \rightarrow aB$ 或 $A \rightarrow a$ ，其中 $A \in V_N$ ， $B \in V_N$ ， $a \in V_T$ 。

例3.1: $\langle \text{标识符} \rangle \rightarrow c \langle \text{字母数字} \rangle$

$\langle \text{字母数字} \rangle \rightarrow c | d | c \langle \text{字母数字} \rangle | d \langle \text{字母数字} \rangle$

$\langle \text{无符号数} \rangle \rightarrow d | d \langle \text{无符号数} \rangle$

$\langle \text{运算符} \rangle \rightarrow + | - | * | / | =$

$\langle \text{界符} \rangle \rightarrow , | ; | (|)$

其中: c 表示字母, d 表示数字

3.3.2 正规式

■ 定义-正规式和它所表示的正规集

设字母表为 Σ ，辅助字母表 $\Sigma' = \{\Phi, \varepsilon, |, ., *, (,)\}$ 。

1. ε 和 Φ 都是 Σ 上的正规式，它们所表示的正规集分别为 $\{\varepsilon\}$ 和 $\{\}$ ；
 2. 任何 $a \in \Sigma$ ， a 是 Σ 上的一个正规式，它所表示的正规集为 $\{a\}$ ；
 3. 假定 e_1 和 e_2 都是 Σ 上的正规式，它们所表示的正规集分别为 $L(e_1)$ 和 $L(e_2)$ ，那么 $(e_1), e_1 | e_2, e_1.e_2, e_1^*$ 也都是正规式，它们所表示的正规集分别为 $L(e_1)$ ， $L(e_1) \cup L(e_2)$ ， $L(e_1)L(e_2)$ 和 $(L(e_1))^*$ 。
 4. 仅由有限次使用上述三步骤而定义的表达式才是 Σ 上的正规式，仅由这些正规式所表示的集合才是 Σ 上的正规集。
- 规定算符的优先顺序为“*”、“.”、“|”。“.”一般可省略不写。“*”、“.”和“|”都是左结合的。

$\{a^n b^m \mid (m, n \geq 1)\}$

$\{a^n b^n \mid (n \geq 1)\}$

■ 例3.2

令 $\Sigma = \{a, b\}$

正规式

a

$a \mid b$

ab

$(a \mid b)(a \mid b)$

a^*

$(a \mid b)^*$

$(a \mid b)^*(aa \mid bb)(a \mid b)^*$

正规集

$\{a\}$

$\{a, b\}$

$\{ab\}$

$\{aa, ab, ba, bb\}$

$\{\varepsilon, a, aa, \dots \text{任意个 } a \text{ 的串}\}$

$\{\varepsilon, a, b, aa, ab, \dots \text{所有由 } a \text{ 和 } b \text{ 组成的串}\}$

$\{\Sigma^* \text{上所有含有两个相继的 } a \text{ 或两个相继的 } b \text{ 组成的串}\}$

■ 又例3.3

1. 令 $\Sigma = \{L, D\}$, 其中 L 代表字母, D 代表数字。
则 Σ 上的正规式

$$r = L(L \mid D)^*$$

定义的正规集为 $\{L, LL, LD, LDD, \dots\}$, 是 “字母打头的字母数字串”, 也就是 Pascal 和多数程序设计语言允许的标识符的词法规则。

2. $\Sigma = \{d, ., e, +, -\}$, 则 Σ 上的正规式
 $d^*(.dd^* \mid \varepsilon)(e(+ \mid - \mid \varepsilon)dd^* \mid \varepsilon)$ 表示的是无符号数的集合。其中 d 为 $0 \sim 9$ 的数字。

- 若两个正规式 e_1 和 e_2 所表示的正规集相同,则说 e_1 和 e_2 等价,写作 $e_1=e_2$ 。
- 设 r,s,t 为正规式, 正规式服从的代数规律有
 1. $r \mid s = s \mid r$ “或” 服从交换律
 2. $r \mid (s \mid t) = (r \mid s) \mid t$ “或” 的可结合律
 3. $(rs)t = r(st)$ “连接” 的可结合律
 4. $r(s \mid t) = rs \mid rt$
 $(s \mid t)r = sr \mid tr$ 分配律
 5. $\varepsilon r = r, r\varepsilon = r$ ε 是“连接”的恒等元素
 6. $r \mid r = r$
 $r^* = \varepsilon \mid r \mid rr \mid \dots$ “或” 的抽取律

例如: $e_1 = a \mid b, e_2 = b \mid a$
 又如: $b(ab)^* = (ba)^*b, (a \mid b)^* = (a^* \mid b^*)^*$

课堂练习

设计表示如下语言的正规表达式：

该语言中的每个字符串由交替的 0 和 1 构成。

$$(01)^* \mid (10)^* \mid 0(10)^* \mid 1(01)^*$$
$$(\varepsilon \mid 1)(01)^*(\varepsilon \mid 0)$$
$$(\varepsilon \mid 0)(10)^*(\varepsilon \mid 1)$$

3.3.3 正规文法到正规式

- 一个正规语言可以由正规文法定义，也可以由正规式定义，对任意一个正规文法，存在一个定义同一个语言的正规式；反之，对每个正规式，存在一个生成同一个语言的正规文法。

■ 对 Σ 上的正规式 r , 存在一个正规 $G=(V_N, V_T, P, S)$ 使 $L(G)=L(r)$ 。

首先, $V_T = \Sigma$, $S \in V_N$, 生成产生式:

$S \rightarrow r$, 且 S 为 G 的开始符。

R.1: 对形如 $A \rightarrow r_1 r_2$ 的产生式, 生成产生式:

$A \rightarrow r_1 B$, $B \rightarrow r_2$ $B \in V_N$

R.2: 对形如 $A \rightarrow r^* r_1$ 的产生式, 生成产生式:

$A \rightarrow rB$, $A \rightarrow r_1$, $B \rightarrow rB$, $B \rightarrow r_1$ $B \in V_N$

R.3: 对形如 $A \rightarrow r_1 | r_2$ 的产生式, 生成产生式:

$A \rightarrow r_1$, $A \rightarrow r_2$

不断应用上述规则进行变换, 直到每个产生式右端只含一个 V_N 。

例3.4: $\Sigma = \{a, d\}$, $r = a(a | d)^*$

正规式到正规文法

正规式	正规文法
$r_1 r_2$	$A \rightarrow r_1 B, B \rightarrow r_2$
$r^* r_1$	$A \rightarrow r B, A \rightarrow r_1, B \rightarrow r B, B \rightarrow r_1$
$r_1 r_2$	$A \rightarrow r_1, A \rightarrow r_2$

例: $\Sigma = \{a, d\}, r = a(a | d)^*$

第1步: $V_T = \{a, d\} \quad S \rightarrow a(a | d)^*$

第2步: $S \rightarrow a(a | d)^*$ 变为:

$S \rightarrow aA, A \rightarrow (a | d)^*$ 利用R.1

第3步: $A \rightarrow (a | d)^*$ 变为:

$A \rightarrow (a | d)B, A \rightarrow \varepsilon, B \rightarrow (a | d)B, B \rightarrow \varepsilon$ 利用R.2

第4步: $A \rightarrow (a | d)B$ 和 $B \rightarrow (a | d)B$ 变为:

$A \rightarrow aB, A \rightarrow dB, B \rightarrow aB, B \rightarrow dB$ 利用R.3

最后得 $G[s]: S \rightarrow aA$

$A \rightarrow aB | dB | \varepsilon$

$B \rightarrow aB | dB | \varepsilon$

正规式 \rightarrow 正规文法的示例

■ 例: $\Sigma = \{a, d\}$, $r = a(a|d)^*$

第1步: $V_T = \{a, d\}$ $S \rightarrow a(a|d)^*$

第2步: $S \rightarrow a(a|d)^*$ 变为:

$S \rightarrow aA$, $A \rightarrow (a|d)^*$ 利用R.1

第3步: $A \rightarrow (a|d)^*$ 变为:

$A \rightarrow (a|d)B$, $A \rightarrow \varepsilon$, $B \rightarrow (a|d)B$, $B \rightarrow \varepsilon$ 利用R.2

第4步: $A \rightarrow (a|d)B$ 和 $B \rightarrow (a|d)B$ 变为:

$A \rightarrow aB$, $A \rightarrow bB$, $B \rightarrow aB$, $B \rightarrow bB$ 利用R.3

最后得 $G[s]$: $S \rightarrow aA$

$A \rightarrow aB|dB|\varepsilon$

$B \rightarrow aB|dB|\varepsilon$

■ 将正规文法转换成正规式

	文法产生式	正规式
规则1	$A \rightarrow xB \quad B \rightarrow y$	$A=xy$
规则2	$A \rightarrow xA y$	$A=x^*y$
规则3	$A \rightarrow x \quad A \rightarrow y$	$A=x y$

例3.5:文法G[S]:

$S \rightarrow aA$

$S \rightarrow a$

$A \rightarrow aA$

$A \rightarrow dA$

$A \rightarrow a$

$A \rightarrow d$

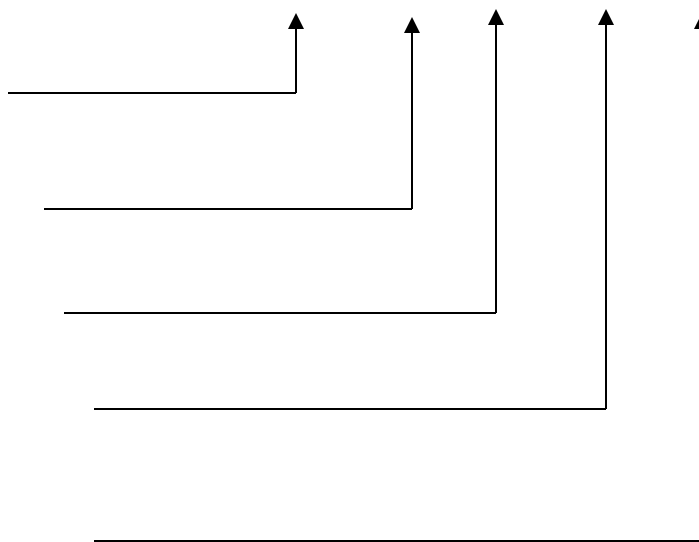
3.4 有穷自动机

- 有穷自动机(也称有限自动机)作为一种识别装置, 它能准确地识别正规集, 即识别正规文法所定义的语言和正规式所表示的集合, 引入有穷自动机这个理论, 正是为词法分析程序的自动构造寻找特殊的方法和工具。
- 有穷自动机分为两类:
 - 确定的有穷自动机(Deterministic Finite Automata)
 - 不确定的有穷自动机(Nondeterministic Finite Automata)

3.4.1 确定的有穷自动机DFA

一个确定有限状态自动机 **DFA**
(*deterministic finite automata*) 是一个五元组
 $A = (K, \Sigma, f, S, Z)$.

- 有限状态集
- 有限输入符号集
- 转移函数
- 一个开始状态
- 一个终态集合



$f : K \times \Sigma \rightarrow K$
 $S \in K$
 $Z \subseteq K$

例3.6

DFA $M = (\{S, U, V, Q\}, \{a, b\}, f, S, \{Q\})$, 其中 f 定义为:

$f(S, a) = U$

$f(V, a) = U$

$f(S, b) = V$

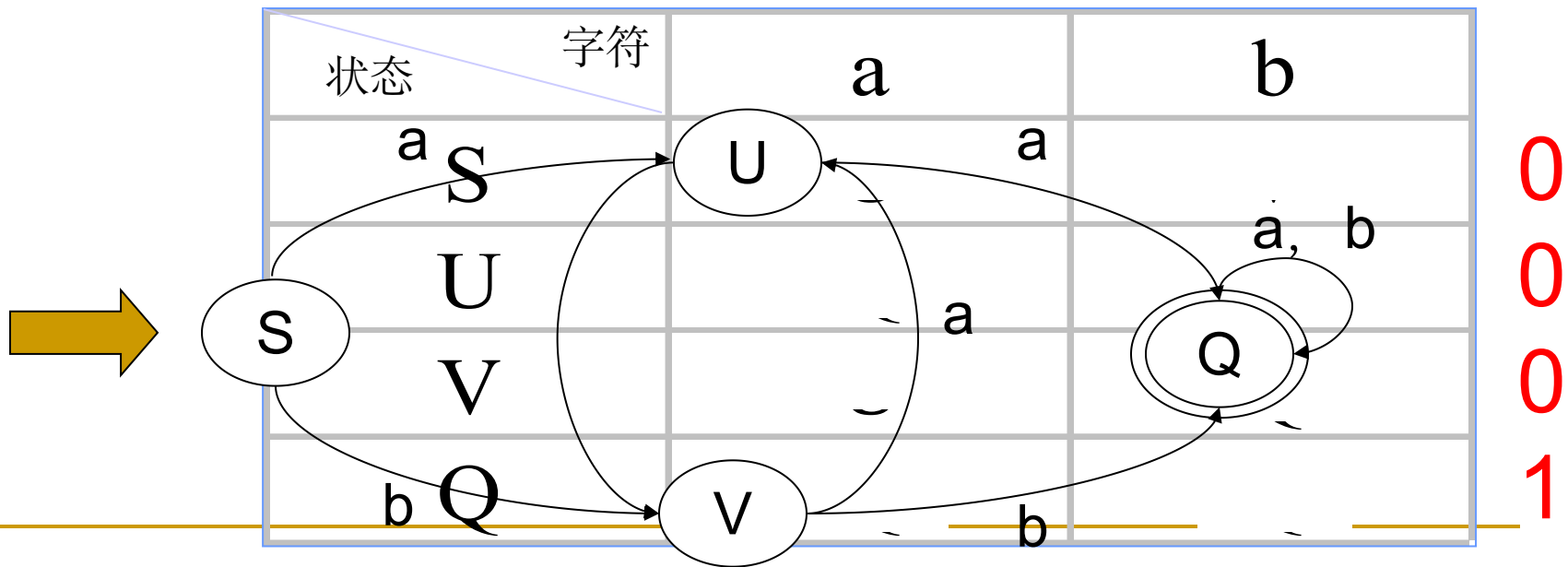
$f(V, b) = Q$

$f(U, a) = Q$

$f(Q, a) = Q$

$f(U, b) = V$

$f(Q, b) = Q$



字符串t能被M接受

- 一个输入字符串 t (将它表示成 t_1t_x 的形式, 其中 $t_1 \in \Sigma$, $t_x \in \Sigma^*$)在DFA M 上运行的定义为:
 - $f(Q, t_1t_x) = f(f(Q, t_1), t_x)$ 其中 $Q \in K$
 - 扩充转换函数 f 是 $K \times \Sigma^* \rightarrow K$ 上的映射, 且: $f(k_i, \varepsilon) = k_i$
- Σ^* 上的字符串 t 被 M 接受
 - 若 $t \in \Sigma^*$, $f(S, t) = P$, 其中 S 为 M 的开始状态, $P \in Z$, Z 为终态集。
 - 则称 t 为DFA M 所接受(识别)。

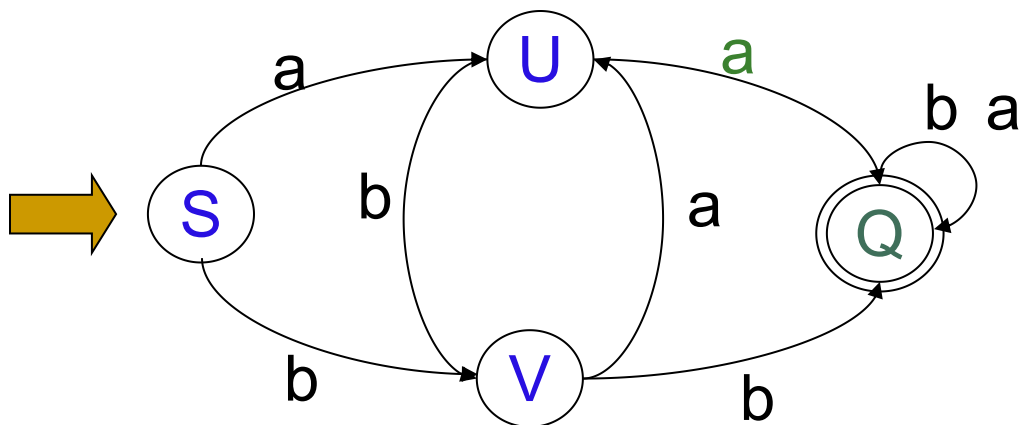
例：证明 $t=baab$ 被下图的DFA所接受。

$$\begin{aligned} f(S, baab) &= f(f(S, b), aab) \\ &= f(V, aab) = f(f(V, a), ab) \\ &= f(U, ab) = f(f(U, a), b) \\ &= f(Q, b) = Q \end{aligned}$$

Q属于终态。

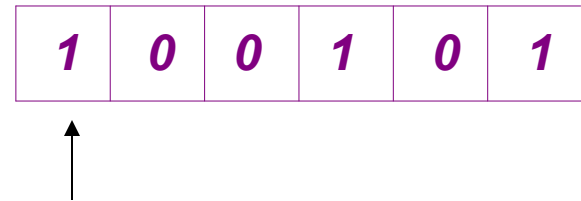
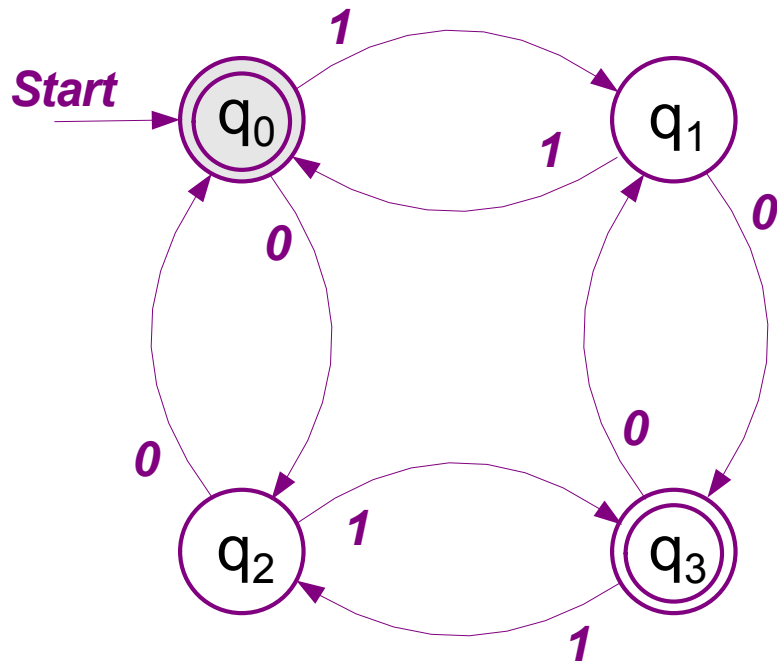
得证。

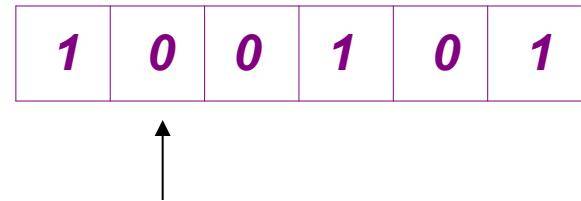
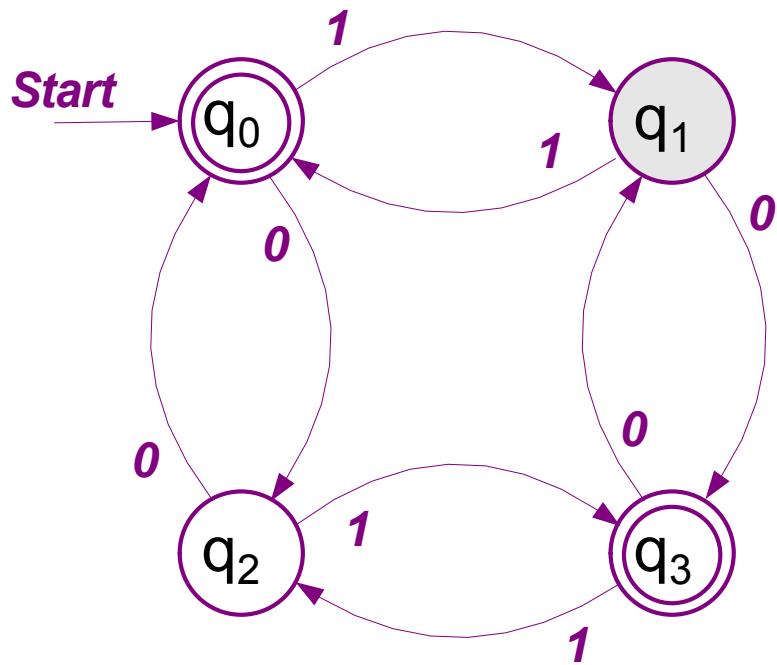
问 $t=ababab$ 能被下图的DFA所接受？

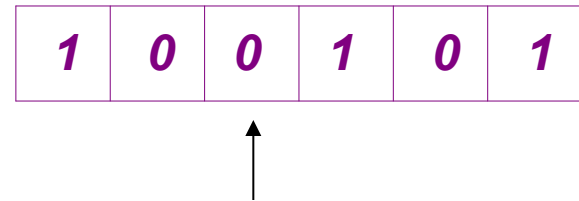
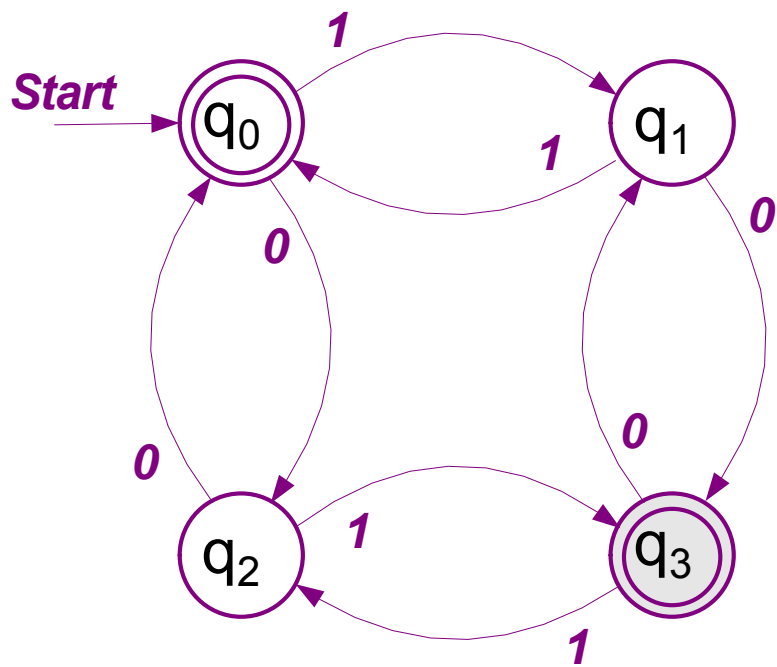


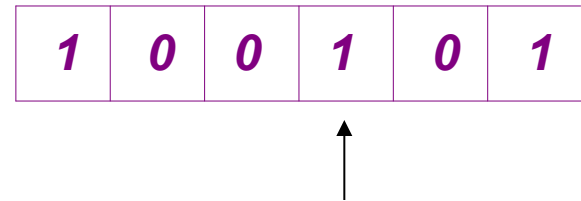
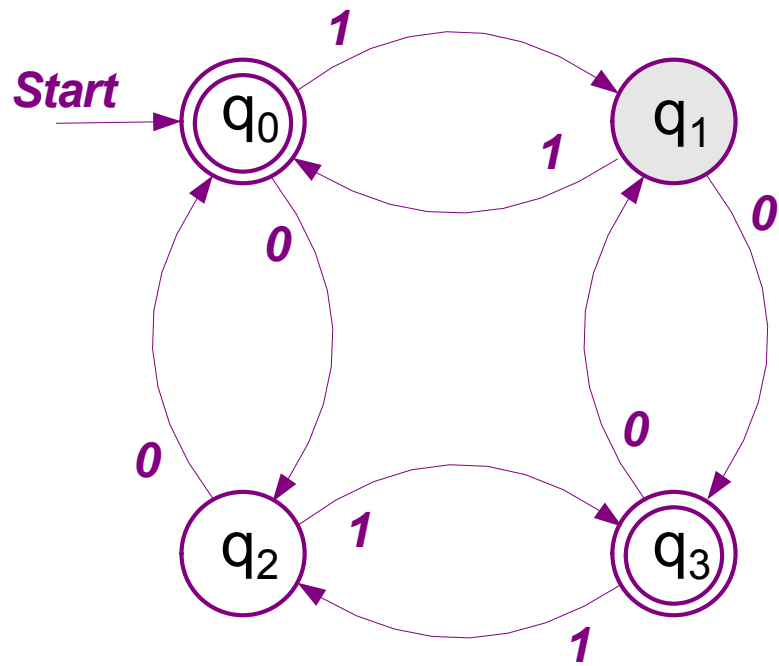
- **DFA $M = (K, \Sigma, f, S, Z)$ 的行为的模拟程序**
 $K := S;$
 $c := \text{getchar}();$
while $c \neq \text{eof}$ do
 $\{K := f(K, c);$
 $c := \text{getchar}();$
 $\};$
if K is in Z then return ('yes')
 else return ('no')

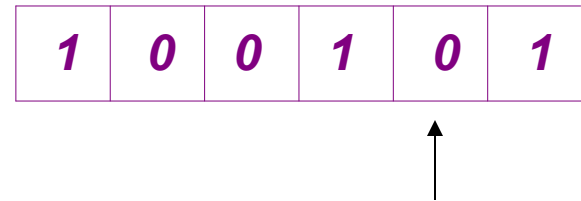
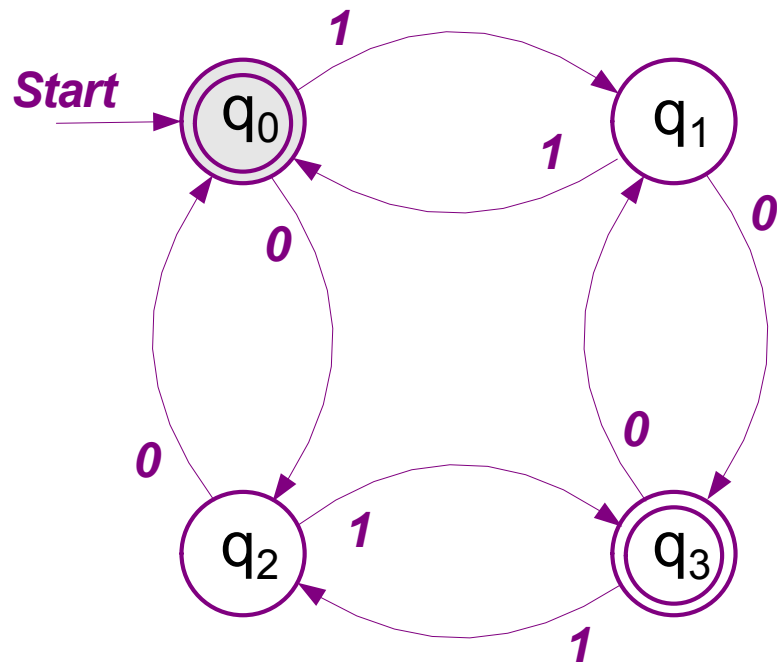
✧ DFA如何接受输入符号串

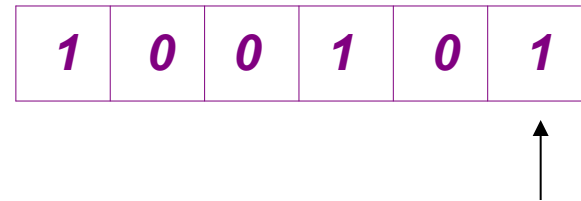
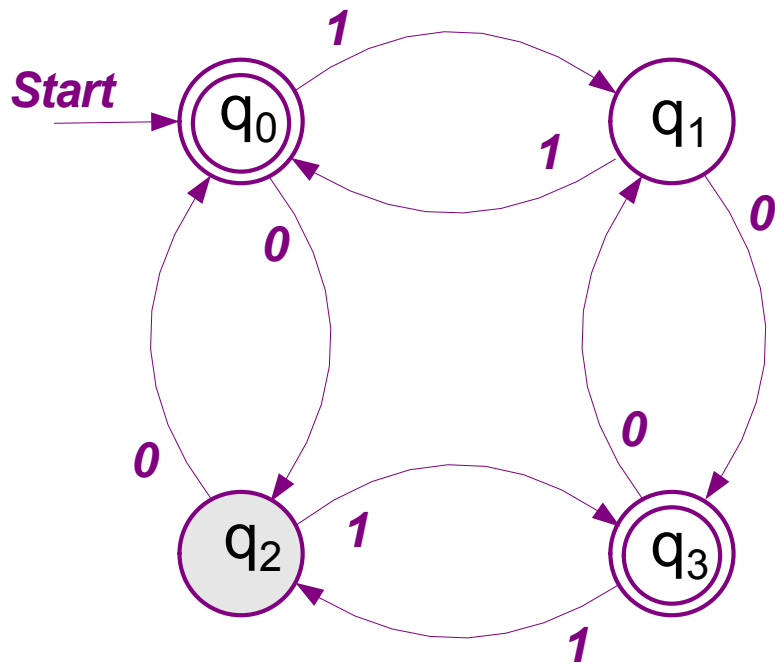


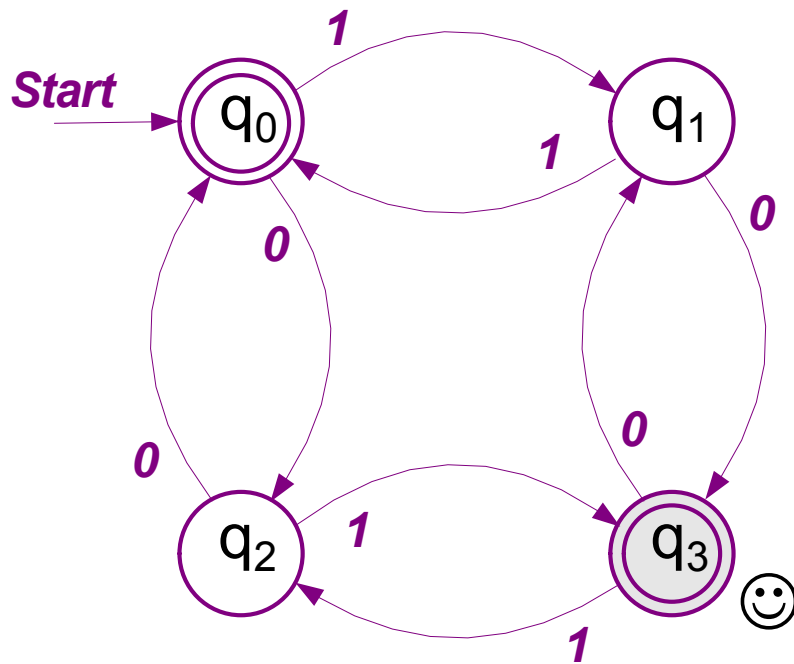


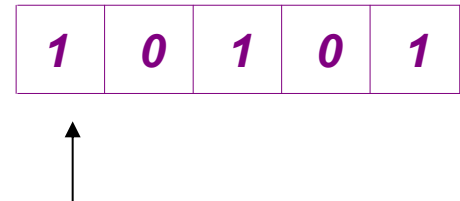
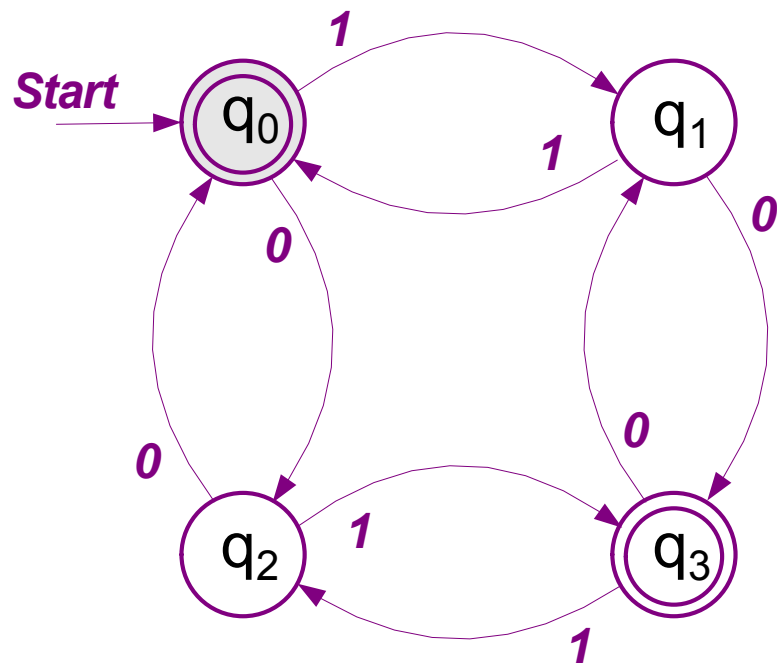


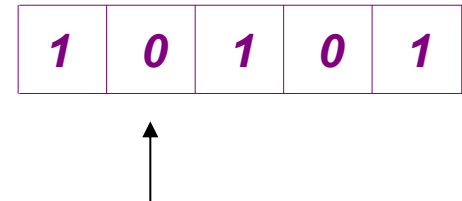
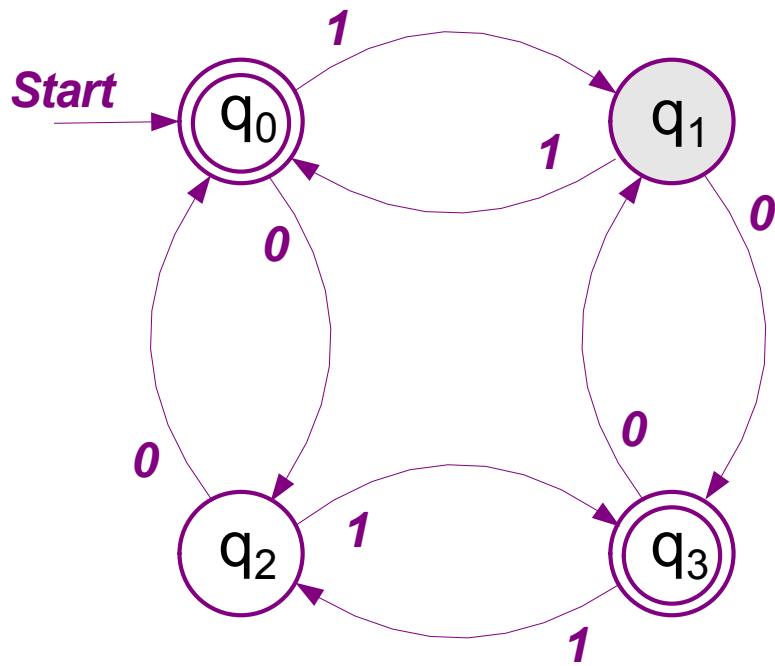


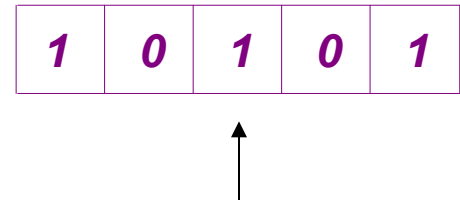
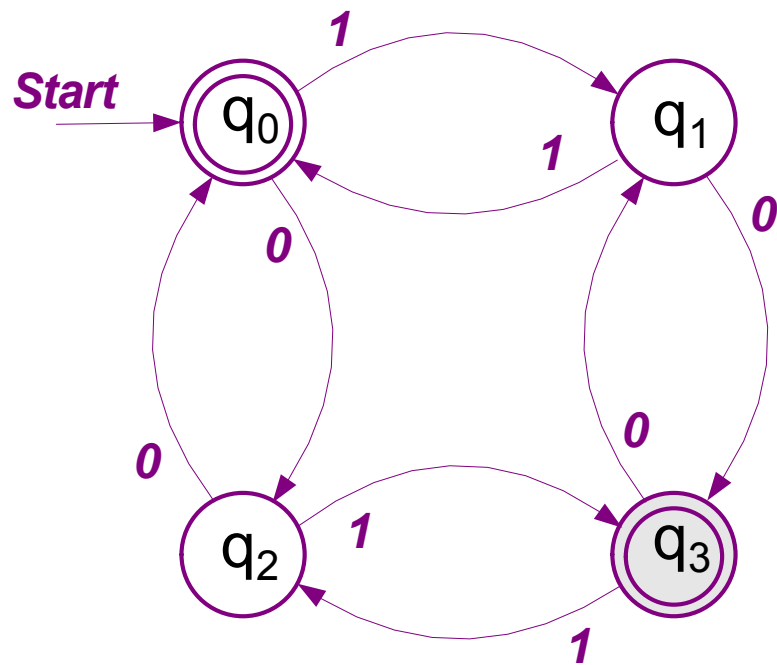


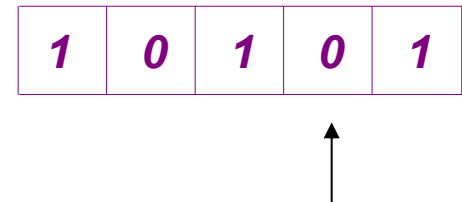
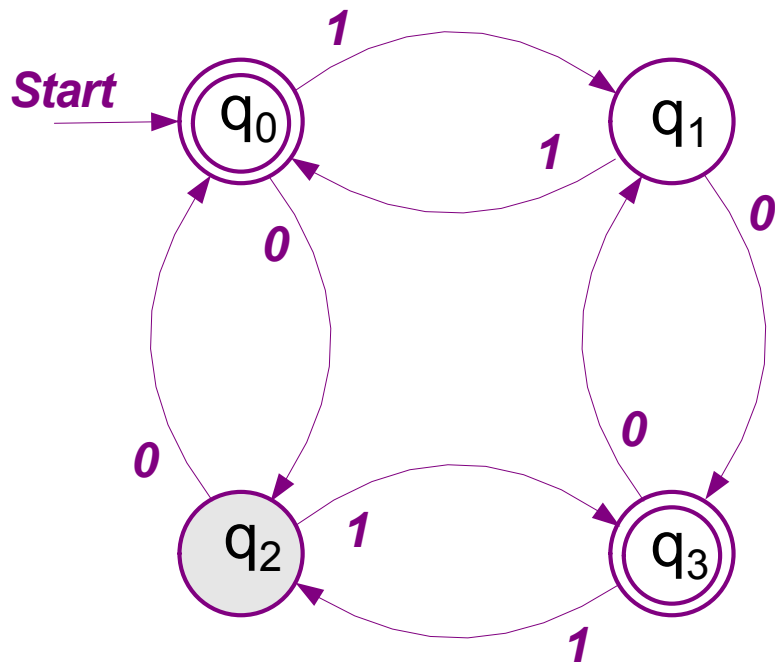


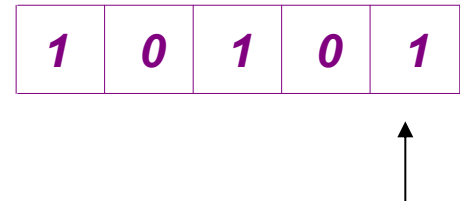
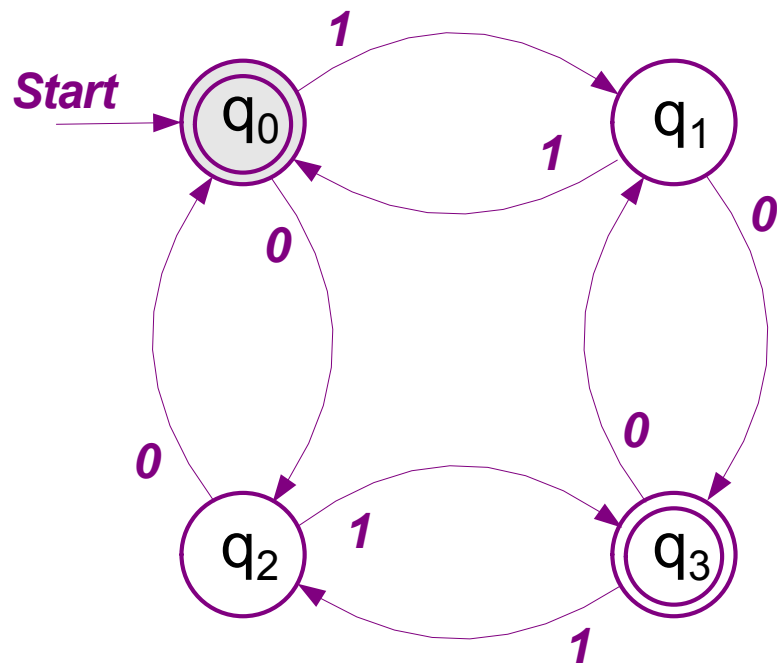


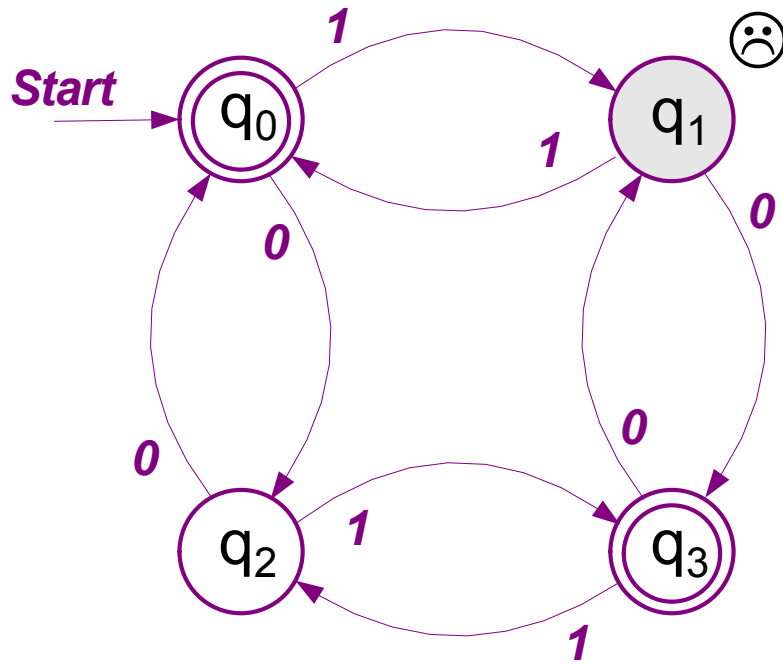












- **DFA M 所能接受的符号串的全体记为 $L(M)$ 。**
对于任何两个有穷自动机 M 和 M' ，如果 $L(M)=L(M')$ ，则称 M 与 M' 是等价的。
- **结论**
 Σ 上一个符号串集 $V \subset \Sigma^*$ 是正规的，当且仅当存在一个 Σ 上的确定有穷自动机 M ，使得 $V=L(M)$ 。

- DFA的确定性表现在
转换函数 $f:K \times \Sigma \rightarrow K$ 是一个单值函数，也就是说，对任何状态 $k \in K$ 和输入符号 $a \in \Sigma$ ， $f(k,a)$ 唯一地确定了下一个状态。从状态转换图来看，若字母表 Σ 含有 n 个输入字符，那么任何一个状态结点最多有 n 条弧射出，而且每条弧以一个不同的输入字符标记。

3.4.2 不确定的有穷自动机NFA

■ 定义

NFA $M = \{K, \Sigma, f, S, Z\}$, 其中

1. K 为状态的有穷非空集,
2. Σ 为有穷输入字母表,
3. f 为 $K \times \Sigma^*$ 到 K 的 **子集** 的一种映射,
4. $S \subseteq K$ 是 **初始状态集**,
5. $Z \subseteq K$ 为终止状态集。

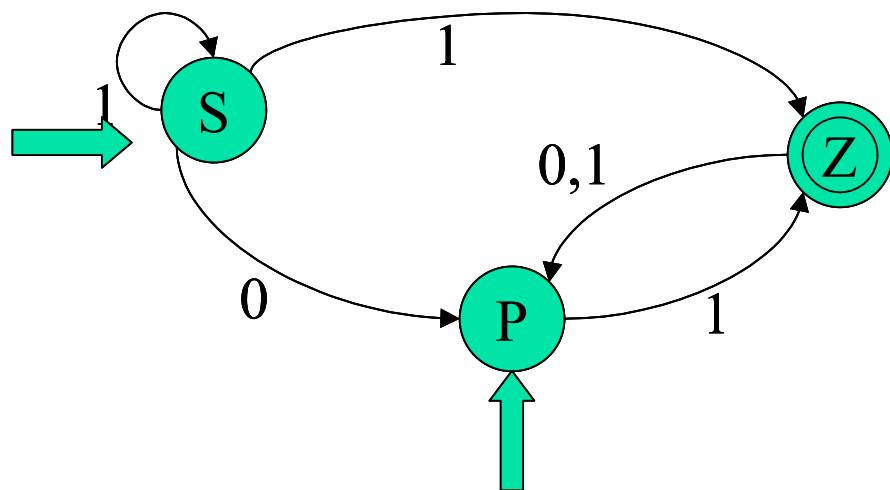
DFA $M = (K, \Sigma, f, S, Z)$

1. K 是一个有穷非空集;
2. Σ 是一个有穷输入字母表
3. f 是转换函数, 是在 $K \times \Sigma \rightarrow K$ 上的映射
4. $S \in K$ 是 **唯一的一个初态**;
5. $Z \subseteq K$ 是一个终态集。

■ 例子

NFA $M = (\{S, P, Z\}, \{0, 1\}, f, \{S, P\}, \{Z\})$, 其中

$f(S, 0) = \{P\}$, $f(Z, 0) = \{P\}$,
 $f(P, 1) = \{Z\}$, $f(Z, 1) = \{P\}$,
 $f(S, 1) = \{S, Z\}$



	0	1	
S	{P}	{S,Z}	0
P	{}	{Z}	0
Z	{P}	{P}	1

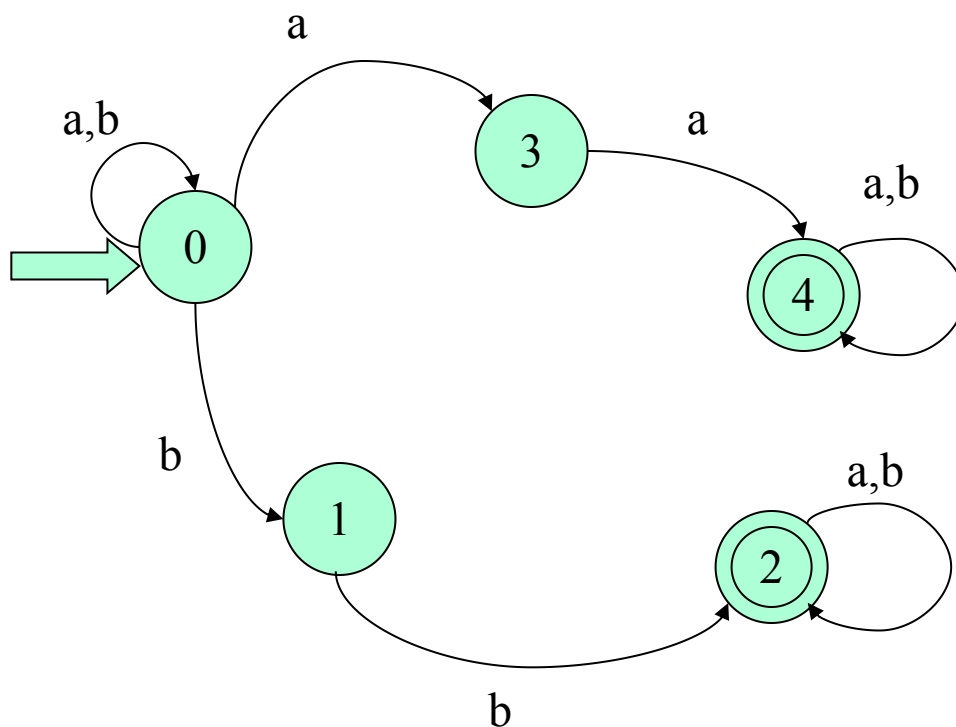
简化为

	0	1	
S	P	S,Z	0
P	.	Z	0
Z	P	P	1

NFA的示例2

■ NFA $M = (\{0,1,2,3,4\}, \{a,b\}, f, \{0\}, \{2,4\})$ 。其中

- $f(0,a) = \{0,3\}$
- $f(0,b) = \{0,1\}$
- $f(1,b) = \{2\}$
- $f(2,a) = \{2\}$
- $f(2,b) = \{2\}$
- $f(3,a) = \{4\}$
- $f(4,a) = \{4\}$
- $f(4,b) = \{4\}$



- Σ^* 上的符号串 t 被NFA M 接受也可以这样理解:
对于 Σ^* 中的任何一个串 t , 若存在一条从某一初态结到某一终态结的道路, 且这条道路上所有弧的标记字依序连接成的串(不理采那些标记为 ϵ 的弧)等于 t , 则称 t 可为NFA M 所识别(读出或接受)。
。
- 若 M 的某些结既是初态结又是终态结, 或者存在一条从某个初态结到某个终态结的道路,其上所有弧的标记均为 ϵ , 那么空字可为 M 所接受。
- **NFA和DFA的不同点:**
 - NFA可以具有 ϵ 转移
 - NFA 中 f 为 $K \times \Sigma^*$ 到 K 的子集 (2^K) 的一种映射。

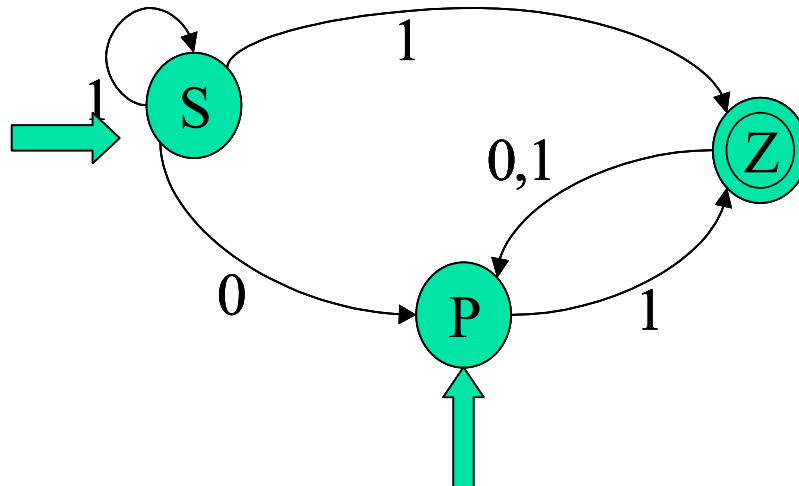
3.4.3 NFA→DFA的转换

- 对每个NFA N 一定存在一个DFA M ，使得 $L(M)=L(N)$ 。对每个NFA N 存在着与之等价的DFA M 。

与某一NFA等价的DFA不唯一。

- 基本思路（子集法）

DFA的每一个状态对应NFA的一组状态。DFA使用它的状态去记录在NFA读入一个输入符号后可能达到的所有状态。



对状态集合I的几个有关运算

- 状态集合I的 ϵ -闭包, 表示为 $\epsilon\text{-closure}(I)$ 。它定义为一状态集, 是状态集I中的任何状态S经任意条 ϵ 弧而能到达的状态的集合。

状态集合I中的任何状态S都属于 $\epsilon\text{-closure}(I)$ 。

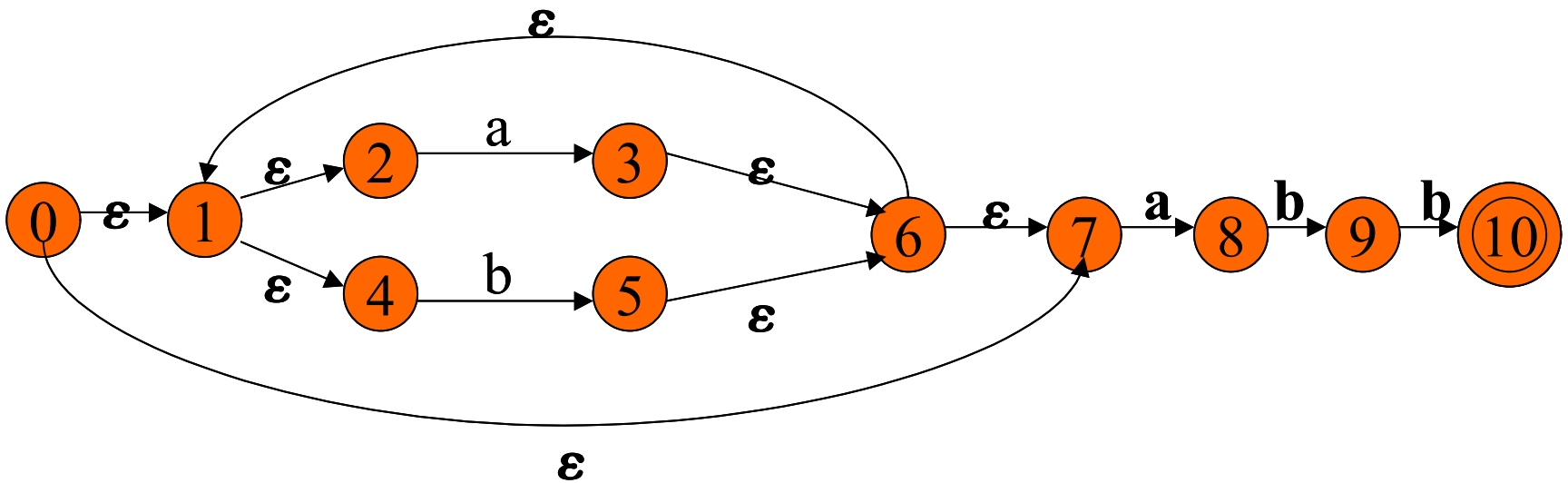
- 状态集合I的a弧转换, 表示为 $\text{move}(I, a)$ 。它定义为状态集合J, 其中J是所有那些可从I的某一状态经过一条a弧而到达的状态的全体。
- 定义 $I_a = \epsilon\text{-closure}(J)$ 。

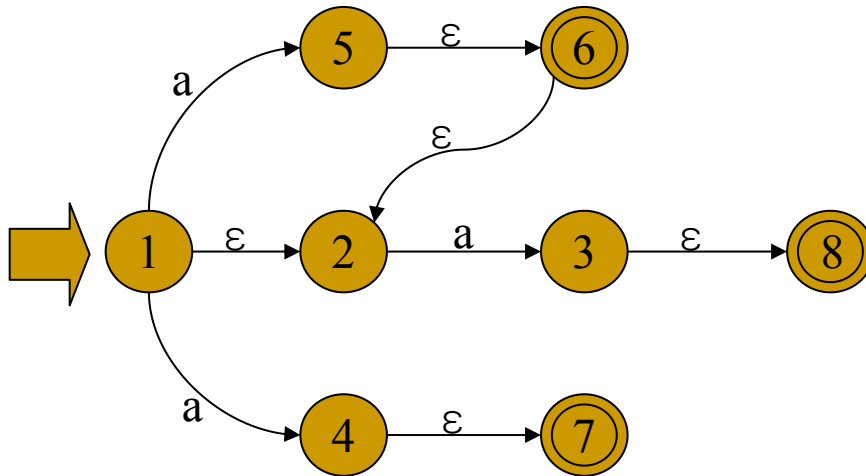
■ 例

ϵ -closure(0)=

■ $\text{move}(\{0,1,2,4,7\},a)=$

■ ϵ -closure($\{3,8\}$)=





1. $I = \{1\}$, ϵ -closure(I) =

2. $I = \{5\}$, ϵ -closure(I) =

3. $I = \{1, 2\}$ move(I, a) =

$I_a =$

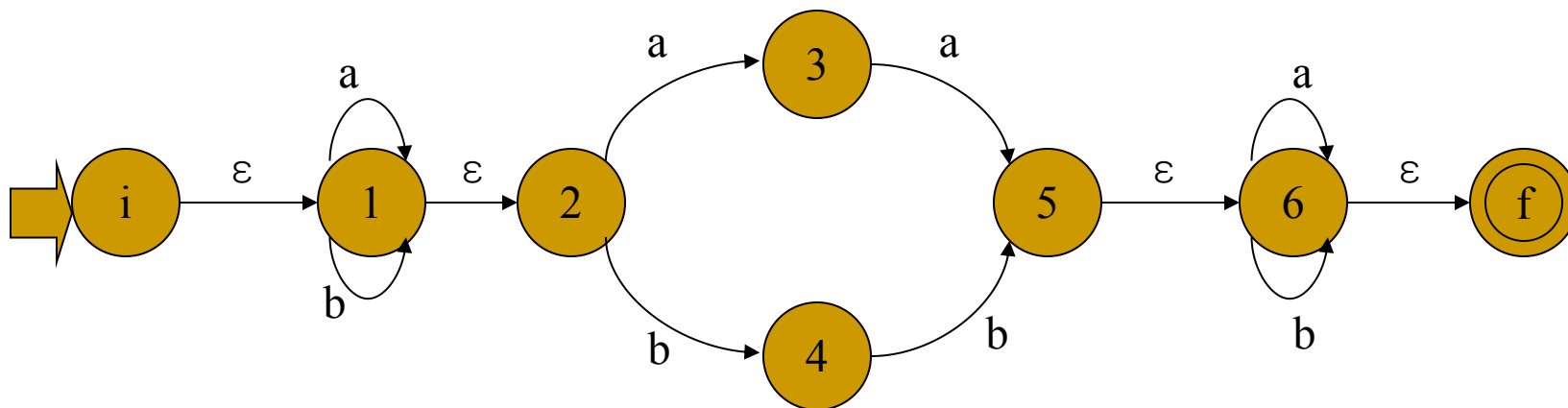
■ NFA确定化算法

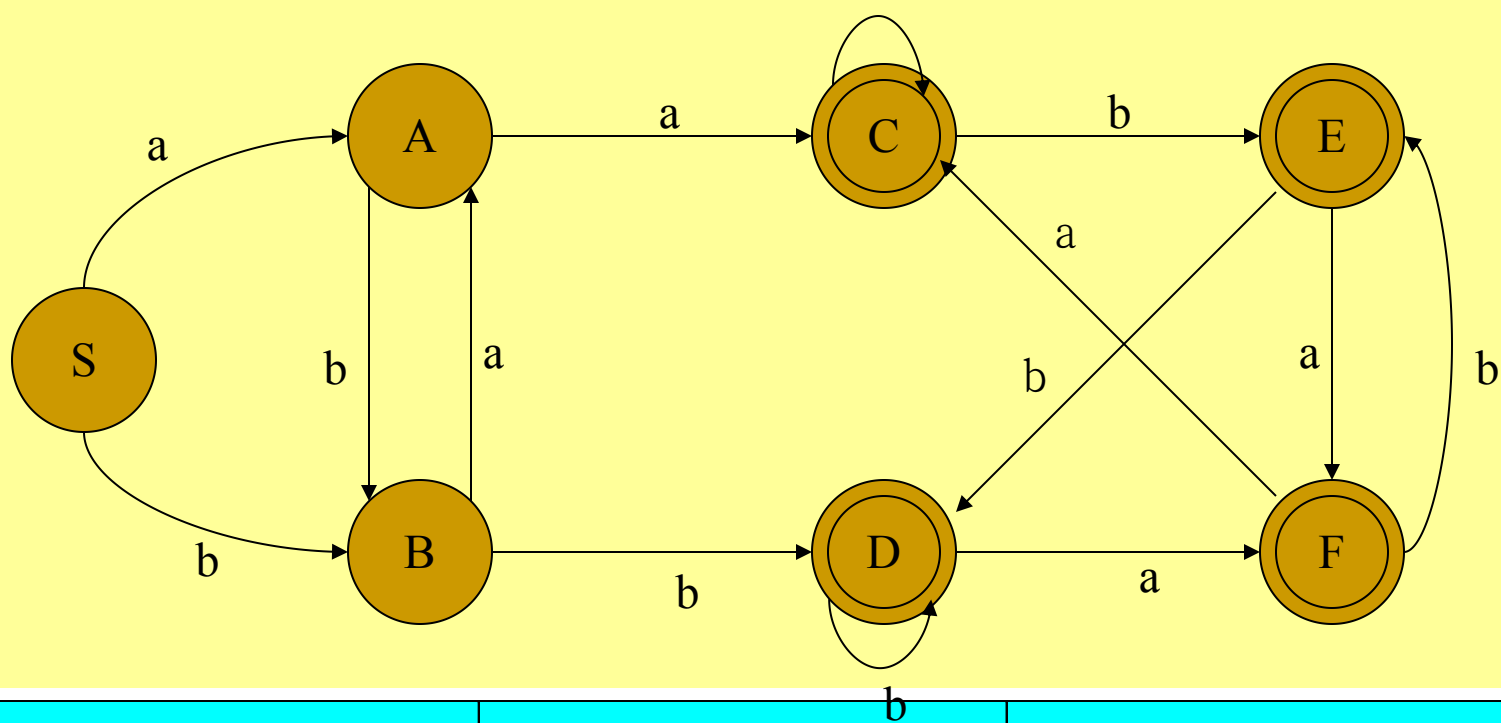
假设NFA $N=(K, \Sigma, f, K_0, K_t)$ 按如下办法构造一个DFA $M=(S, \Sigma, d, S_0, S_t)$, 使得 $L(M)=L(N)$:

1. M 的状态集 S 由 K 的一些子集组成。用 $[S_1, S_2, \dots, S_j]$ 表示 S 的元素, 其中 S_1, S_2, \dots, S_j 是 K 的状态。并且约定, 状态 S_1, S_2, \dots, S_j 是按某种规则排列的, 即对于子集 $\{S_1, S_2\}=\{S_2, S_1\}$ 来说, S 的状态就是 $[S_1 S_2]$;
2. M 和 N 的输入字母表是相同的, 即是 Σ ;
3. 转换函数是这样定义的:
 $d([S_1, S_2, \dots, S_j], a) = [R_1, R_2, \dots, R_t]$, 其中
 $[R_1, R_2, \dots, R_t] = \varepsilon\text{-closure}(\text{move}(\{S_1, S_2, \dots, S_j\}, a))$
4. $S_0 = \varepsilon\text{-closure}(K_0)$ 为 M 的开始状态;
5. $S_t = \{[S_i, S_k, \dots, S_e], \text{ 其中 } [S_i, S_k, \dots, S_e] \in S \text{ 且 } \{S_i, S_k, \dots, S_e\} \cap K_t \neq \emptyset\}$

NFA确定化算法示例

■ 例子



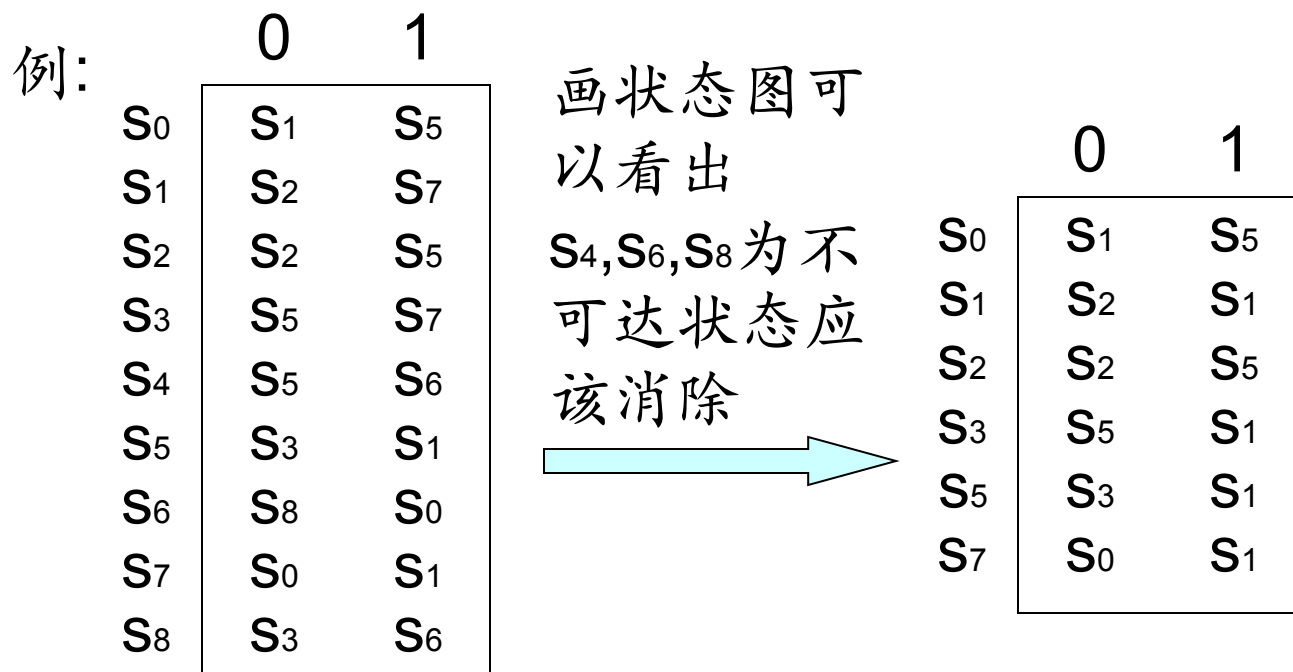


		I_a	I_b
$\{i,1,2\}$	S	$\{1,2,3\}$ A	$\{1,2,4\}$ B
$\{1,2,3\}$	A	$\{1,2,3,5,6,f\}$ C	$\{1,2,4\}$ B
$\{1,2,4\}$	B	$\{1,2,3\}$ A	$\{1,2,4,5,6,f\}$ D
$\{1,2,3,5,6,f\}$	C	$\{1,2,3,5,6,f\}$ C	$\{1,2,4,6,f\}$ E
$\{1,2,4,5,6,f\}$	D	$\{1,2,3,6,f\}$ F	$\{1,2,4,5,6,f\}$ D
$\{1,2,4,6,f\}$	E	$\{1,2,3,6,f\}$ F	$\{1,2,4,5,6,f\}$ D
$\{1,2,3,6,f\}$	F	$\{1,2,3,5,6,f\}$ C	$\{1,2,4,6,f\}$ E

3.4.4 确定有穷自动机的化简

- 对于任一个**DFA**，存在一个唯一的**状态最少的等价的DFA**。
- 一个有穷自动机是化简了的，即是说它没有多余状态并且它的状态中没有两个是互相等价的。
- 一个有穷自动机可以通过消除多余状态和合并等价状态而转换成一个最小的与之等价的有穷自动机。

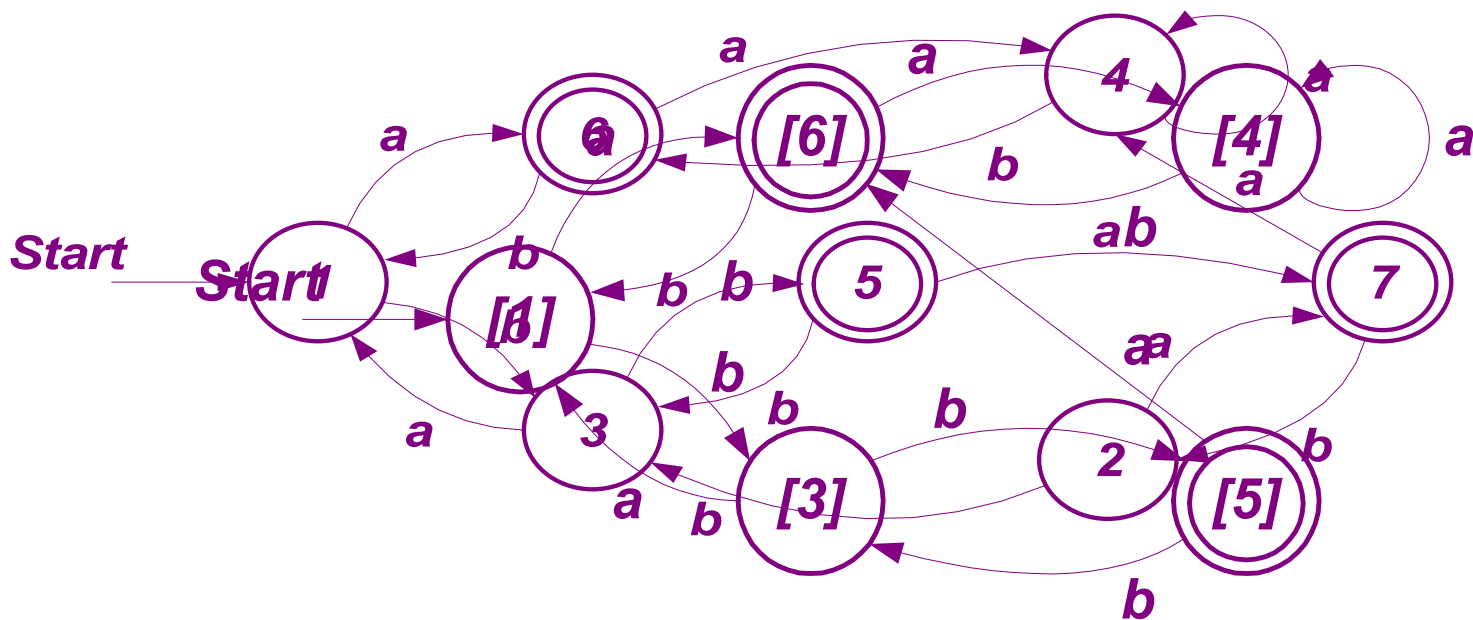
- 所谓有穷自动机的多余状态是指是指这样的状态：从该自动机的开始状态出发，任何输入串也不能到达那个状态。



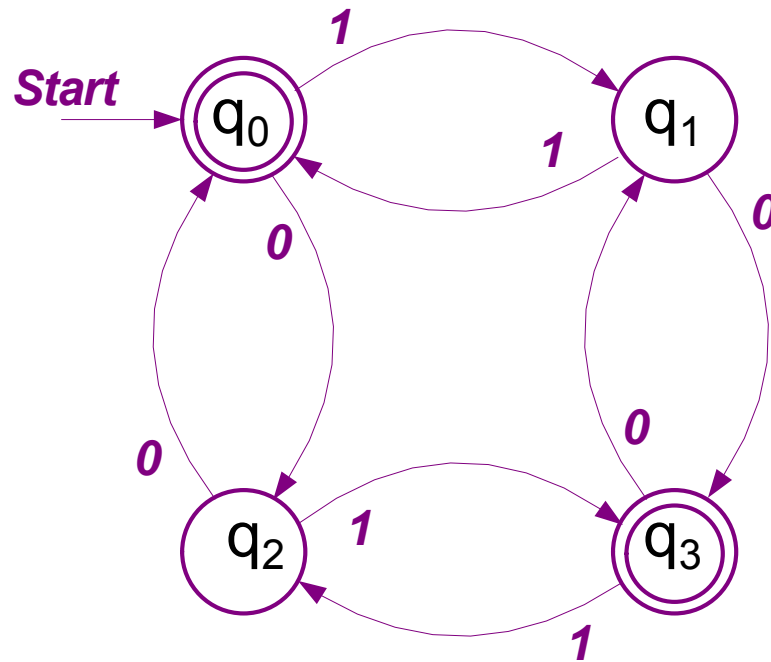
- 状态 s 和 t 的等价条件是:
 - 一致性条件: 状态 s 和 t 必须同时为可接受状态或不接受状态。
 - 蔓延性条件: 对于所有输入符号,状态 s 和 t 必须转换到等价的状态里。
- 有穷自动机的状态 s 和 t 不等价,称这两个状态是可区别的。

■ DFA的最小化算法-“分割法”

把一个**DFA**的状态分成一些不相交的子集，使得任何不同的两子集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。

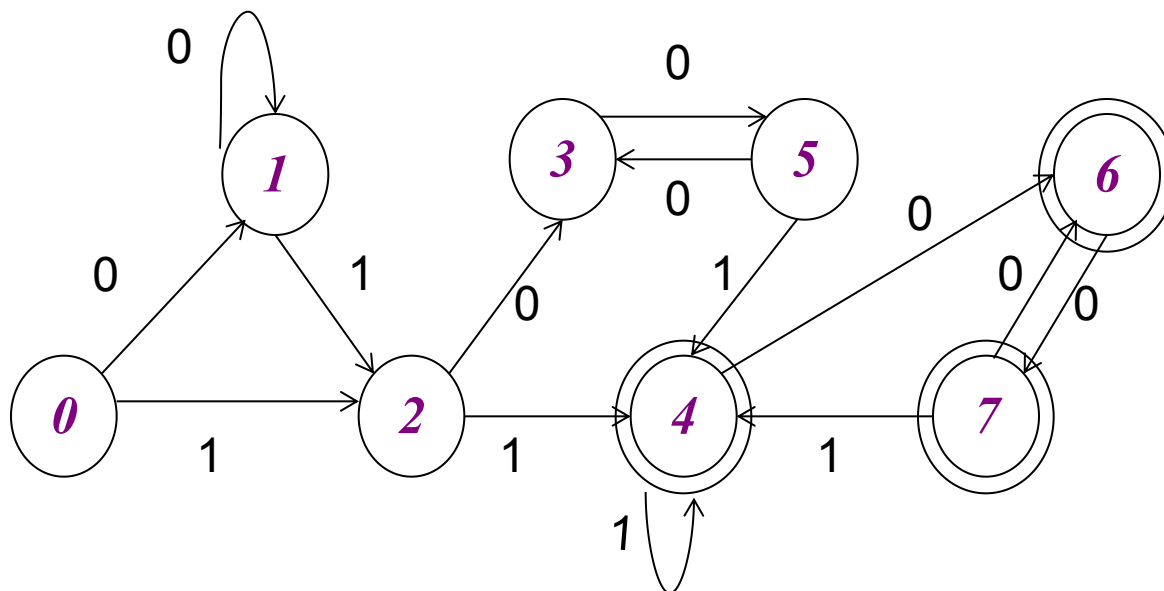


— 最小化下图表示的 DFA



课堂练习

— 最小化下图表示的 DFA



3.5 有穷自动机和正规表达式

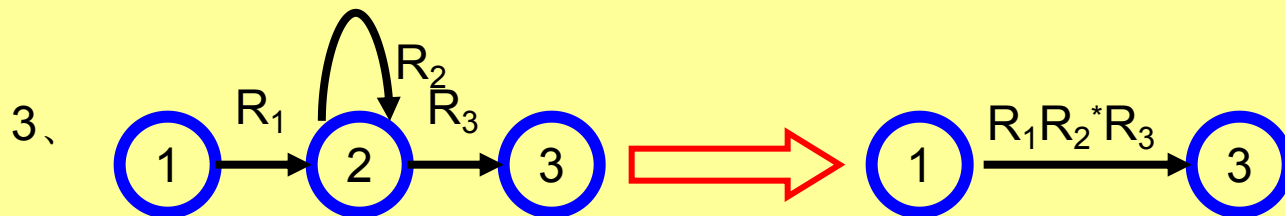
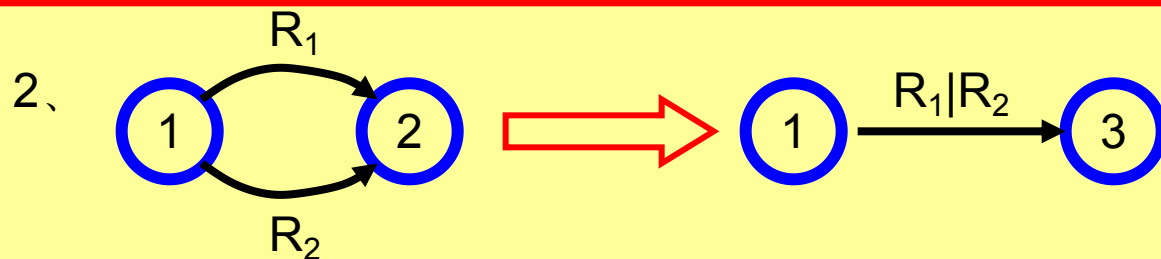
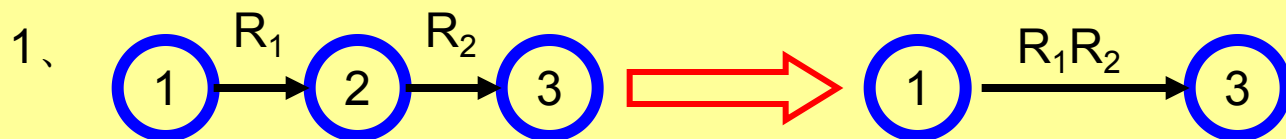
- 有穷自动机和正规表达式的等价性:

1. 对于 Σ 上的一个 **NFA** M , 可以构造一个 Σ 上的正规式 R , 使得 $L(R)=L(M)$ 。

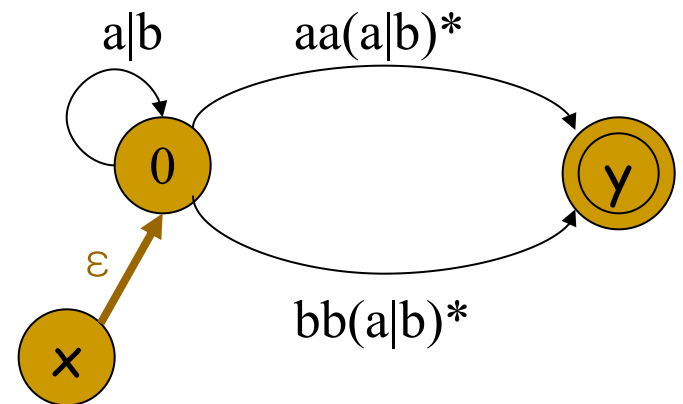
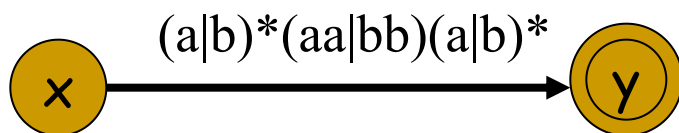
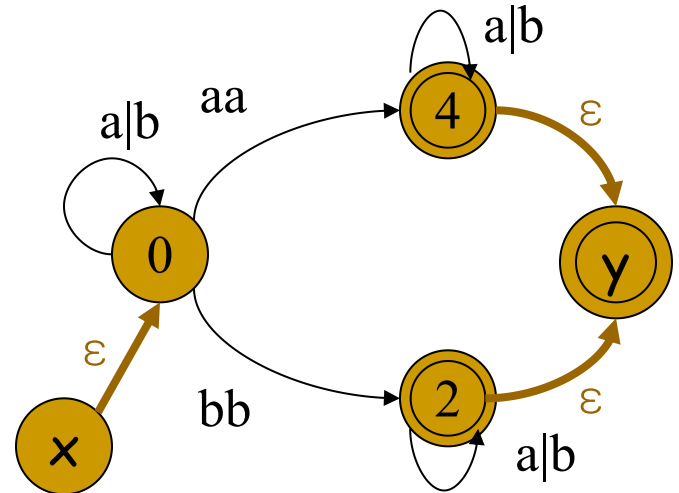
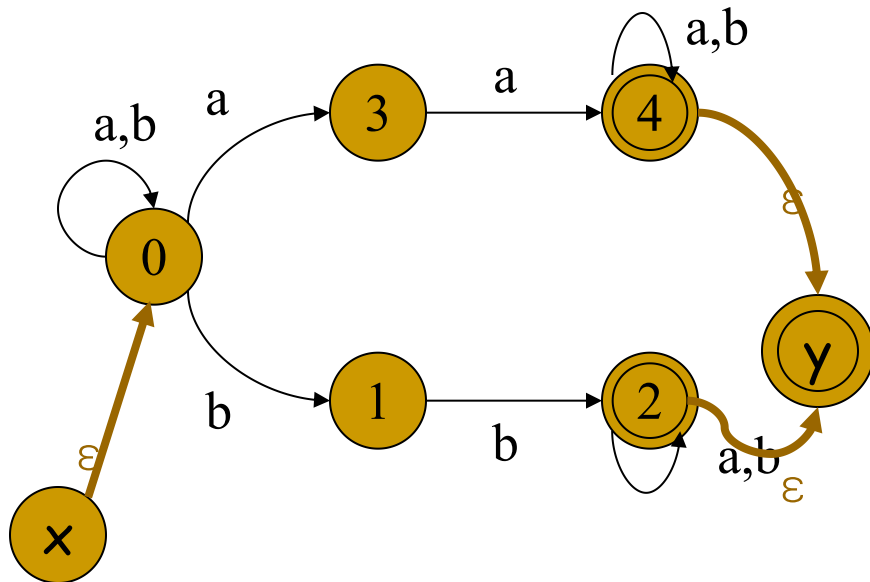
2. 对于 Σ 上的一个正规式 R , 可以构造一个 Σ 上的 **NFA** M , 使得 $L(M)=L(R)$ 。

一、NFA转化为正规式

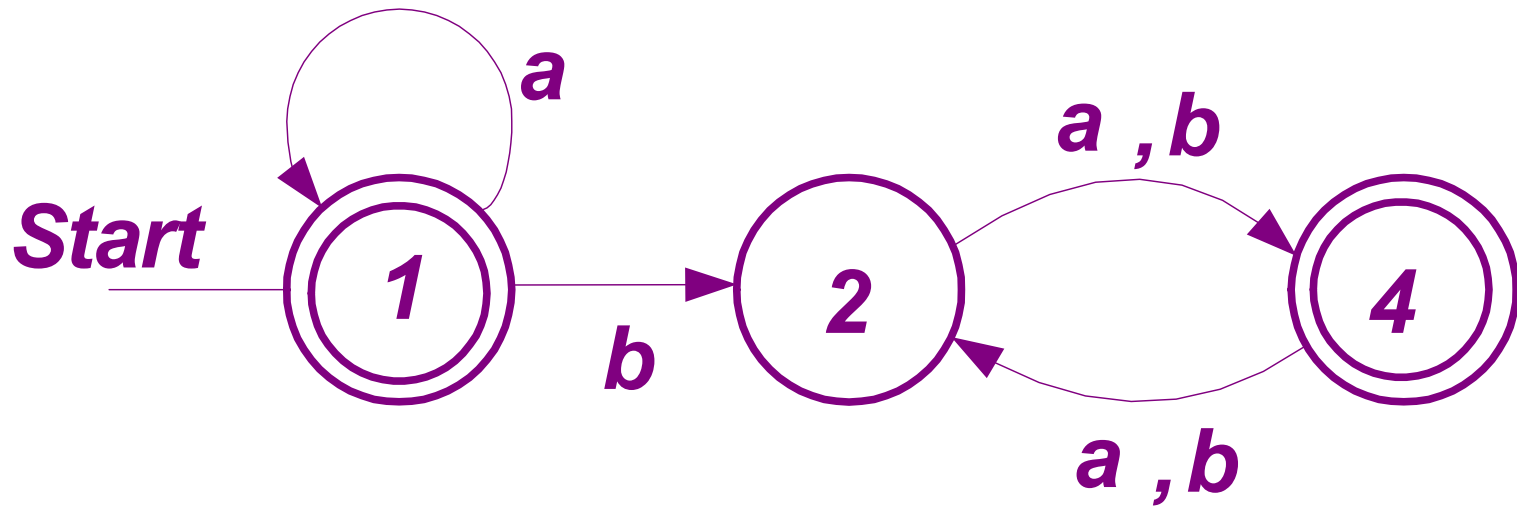
- 第一步：在NFA M 上加两个新的状态 x 和 y ，从 x 出发用 ϵ 弧连接到 M 的所有初态，从 M 的所有终态用 ϵ 弧连接到 y 。这样就形成了一个和 M 等价的 M'



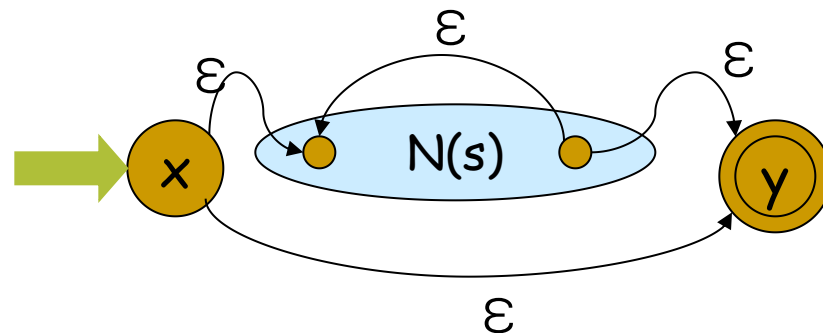
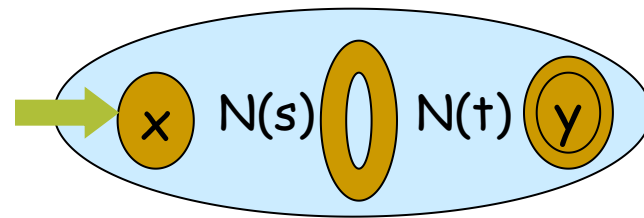
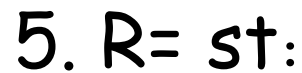
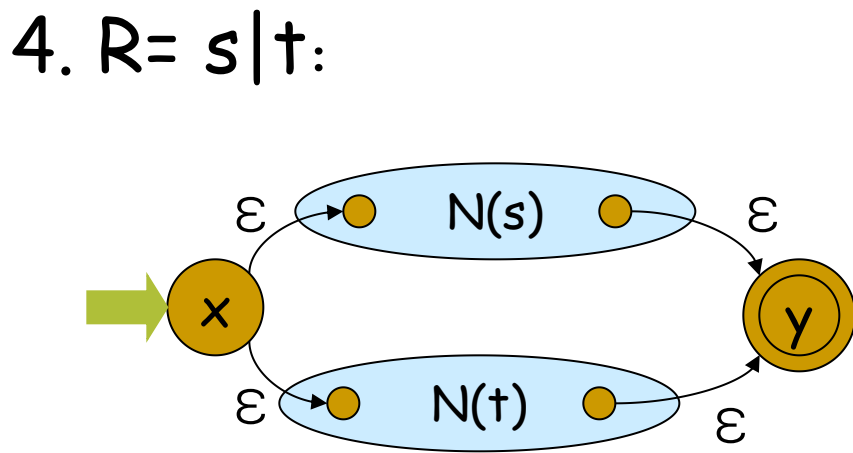
NFA转化为正规式的示例



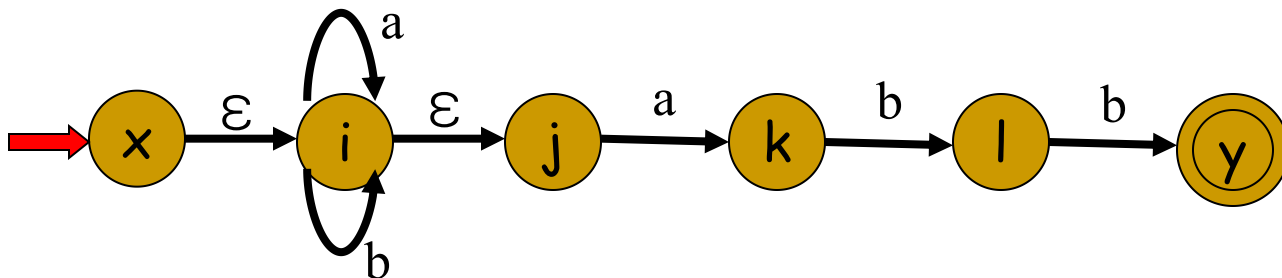
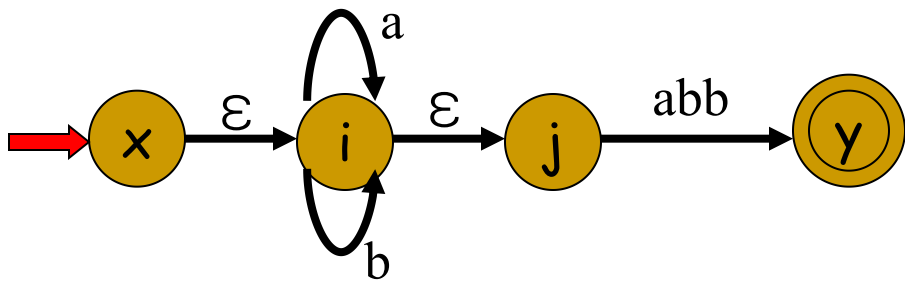
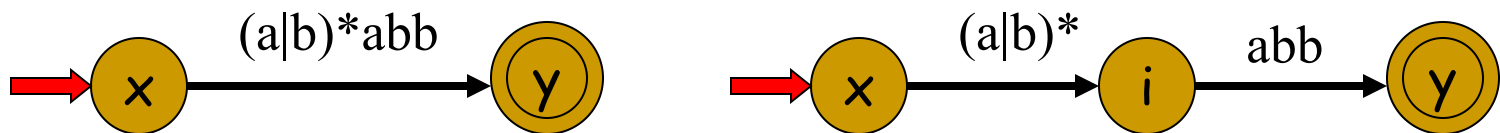
课堂练习



二、正规式转化为NFA



正规式转化为NFA的示例



3.6 正规文法与有穷自动机间的转换

■ 正规文法到NFA

1. 字母表与**G**的终结符相同;
2. 为**G**中的每个非终结符生成**M**的一个状态, **G**的开始符号**S**是开始状态**S**;
3. 增加一个新状态**Z**, 作为**NFA**的终态;
4. 对**G**中的形如**A**→**tB**, 其中**t**为终结符或 ϵ , **A**和**B**为非终结符的产生式, 构造**M**的一个转换函数 $f(A, t) = B$;
5. 对**G**中的形如**A**→**t**的产生式, 构造**M**的一个转换函数 $f(A, t) = Z$ 。

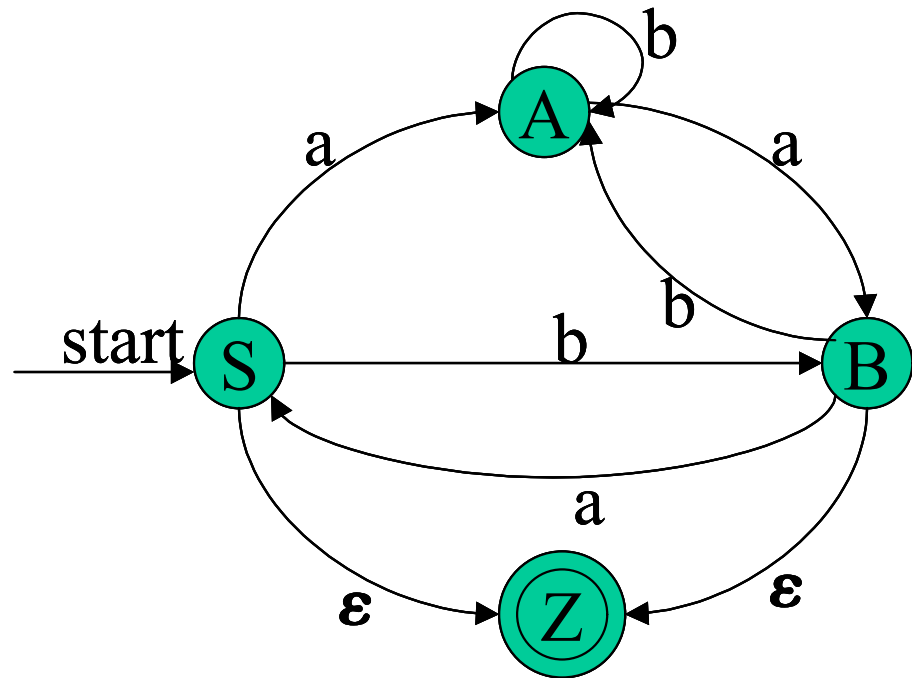
■ 例:求与文法**G[S]**等价的**NFA**

G[S]:

$S \rightarrow aA | bB | \epsilon$

$A \rightarrow aB | bA$

$B \rightarrow aS | bA | \epsilon$



■ NFA到正规文法

1.对转换函数 $f(A,t)=B$ ，可写成一个产生式：

$A \rightarrow tB$

2.对可接受状态 Z ，增加一个产生式： $Z \rightarrow \epsilon$

3.有穷自动机的初态对应于文法的开始符号，
有穷自动机的字母表为文法的终结符号集。

- 例:给出如图**NFA**等价的正则文法**G**

G[A]:

A \rightarrow **aB**

A \rightarrow **bD**

B \rightarrow **bC**

C \rightarrow **aA**

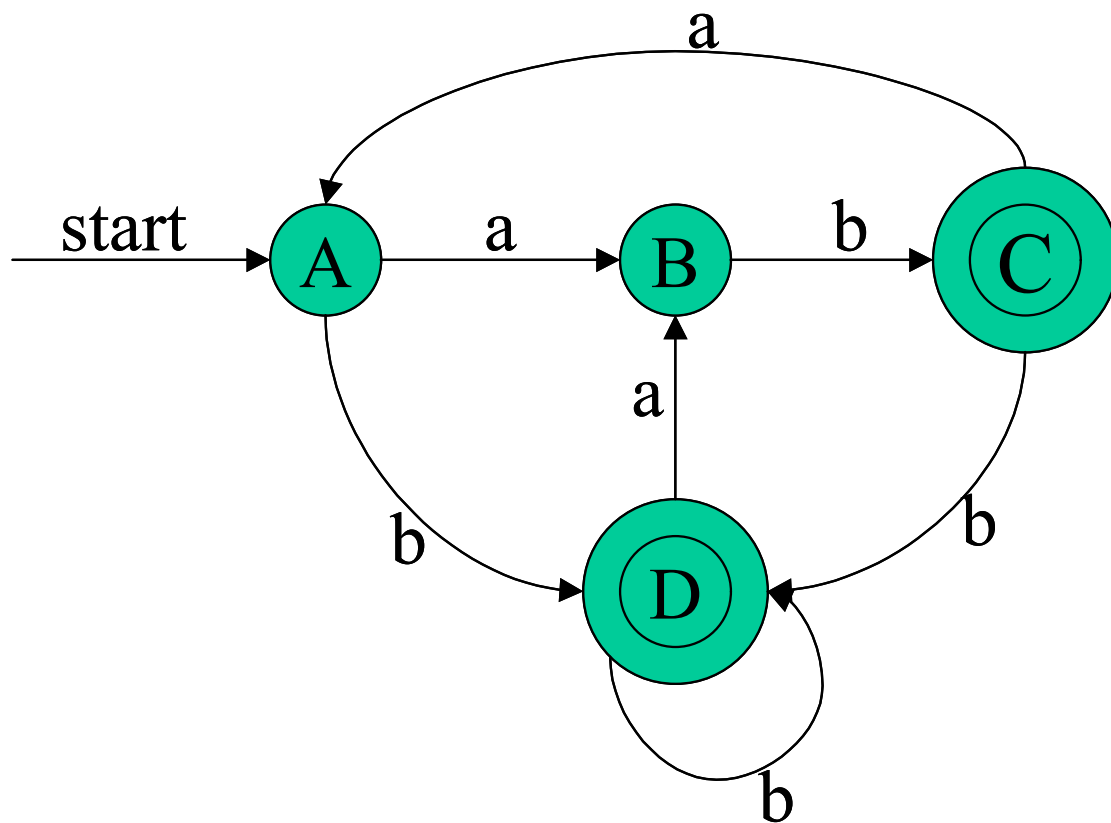
C \rightarrow **bD**

C \rightarrow ϵ

D \rightarrow **aB**

D \rightarrow **bD**

D \rightarrow ϵ

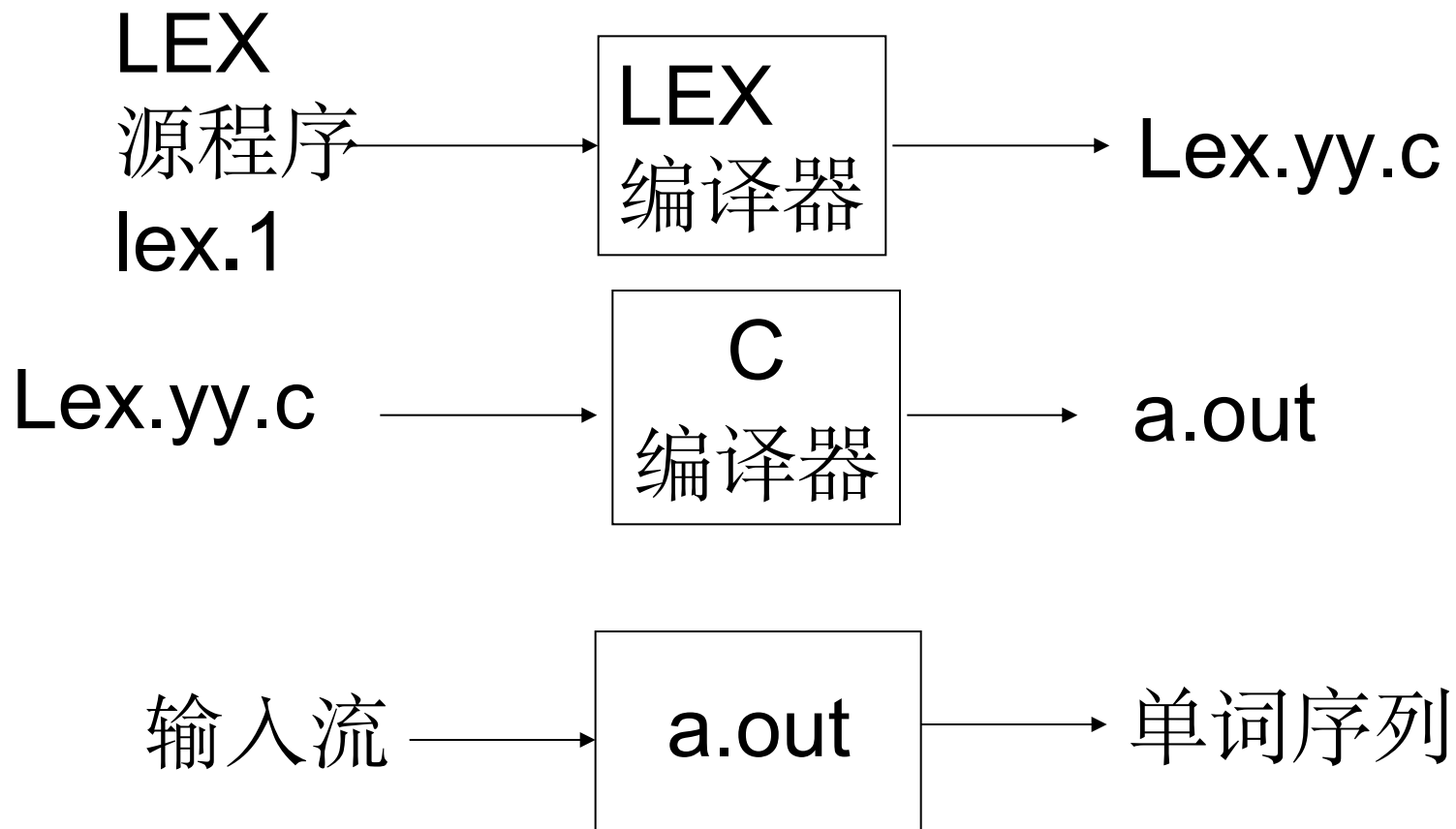


课堂练习

构造正规式 $1(0|1)^*101$ 相应的DFA。

3.7 词法分析程序的自动构造

■ 词法分析程序的自动构造工具LEX简介



课堂练习

构造如下语言的上下文无关文法：

$$(1) \{a^n b^{2n} c^m \mid n, m \geq 0\}$$

$$(2) \{a^n b^{n+m} c^m \mid n, m \geq 0\}$$

给出下列每个正规语言的一个正规表达式

$$1. \{a^n b^m \mid n, m \geq 0 \text{ 且 } n + m \text{ 为偶数} \}$$

$$2. \{w \in \{a, b\}^* \mid w \text{ 中既不包含子串 } aa, \text{ 也不包含子串 } bb \}$$