

---

## 第三章 词法分析

---

# 主要内容

---

本章将讨论词法分析程序的设计原则，单词的描述技术，识别机制及词法分析程序的自动构造原理。

- 3.1 词法分析程序的设计
- 3.2 单词的描述工具
- 3.3 有穷自动机
- 3.4 正规式和有穷自动机的等价性
- 3.5 正规文法和有穷自动机的等价性
- 3.6 词法分析程序的自动构造

## 3.1 词法分析程序的设计

---

- 实现词法分析（lexical analysis）的程序
  - 逐个读入源程序字符并按照构词规则切分成一系列单词。
  - 单词是语言中具有独立意义的最小单位，包括保留字、标识符、运算符、标点符号和常量等。
  - 词法分析是编译过程中的一个阶段，在语法分析前进行。也可以和语法分析结合在一起作为一遍，由语法分析程序调用词法分析程序来获得当前单词供语法分析使用。

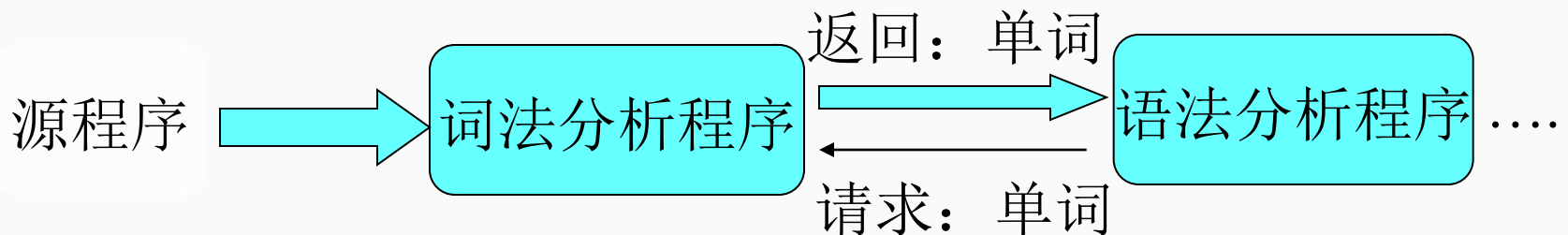
# 词法分析的任务

---

- 词法分析程序的主要任务：
  - 读源程序，产生单词符号
- 词法分析程序的其他任务：
  - 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， .....

# 词法分析程序与语法分析程序的接口方式

1. 词法分析程序作为**独立的一遍**，把字符流的源程序变为单词序列，输出在一个中间文件上，这个文件作为语法分析程序的输入继续编译过程。
2. 把词法分析程序**设计成一个子程序**，每当语法分析程序需要一个单词时，则调用该子程序。词法分析程序每得到一次调用，便从源程序文件读入一些字符，直到识别出一个单词，或说直到下一个单词的第一个字符为止。



# 词法分析的输出

---

- 词法分析的功能是读入源程序，输出单词符号。
- 输出的单词符号一般分成5种：
  1. 关键字(基本字)
  2. 标识符
  3. 常数
  4. 运算符
  5. 界符
- 词法分析程序所输出的单词符号采用二元式表示：(单词种别, 单词自身的值)

# 词法分析的输出举例

设标识符编码为1,

常数为2,

关键字为3,

运算符为4,

界符为5

输出

(3, 'if')

(1, 指向i的符号表入口)

(4, '=')

(2, '5')

(3, 'then')

(1, 指向x的符号表入口)

(4, ':=')

(1, 指向y的符号表入口)

(5, ';')

程序段:

if i=5

then x:=y;

# 词法分析工作从语法分析工作独立出来的原因

---

- 简化设计
- 改进编译效率
- 增加编译系统的可移植性



## 3.3 单词的形式化描述工具

---

程序设计语言中的单词是基本语法符号。单词符号的语法可以用有效的工具加以描述，并且基于这类描述工具，可以建立词法分析技术，进而可以建立词法分析的自动构造方法。

1. 正规文法
2. 正规式
3. 有穷自动机

# 正规文法

---

回顾一下3型文法，设文法 $G=(V_N, V_T, P, S)$ ，若 $P$ 中每一个产生式都是 $A \rightarrow aB$ 或 $A \rightarrow a$ ，其中 $A, B \in V_N$ ， $a \in V_T^*$ ，则 $G$ 是一个3型文法或正规文法。

# 各类单词的正规文法描述

$\langle \text{标识符} \rangle \rightarrow l \mid l \langle \text{字母数字} \rangle$

$\langle \text{字母数字} \rangle \rightarrow l \mid d \mid l \langle \text{字母数字} \rangle \mid d \langle \text{字母数字} \rangle$

$\langle \text{无符号整数} \rangle \rightarrow d \mid d \langle \text{无符号整数} \rangle$

$\langle \text{运算符} \rangle \rightarrow + \mid - \mid * \mid / \mid = \mid \langle \langle \text{等号} \rangle \mid \rangle \langle \text{等号} \rangle \dots$

$\langle \text{等号} \rangle \rightarrow =$

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

其中l表示a~z中的任何一个英文字母，d表示0~9中的任何一个数字

# 各类单词的正规文法描述

- 无符号数

$\langle \text{无符号数} \rangle \rightarrow d \langle \text{余留无符号数} \rangle \mid . \langle \text{十进小数} \rangle \mid e \langle \text{指数部分} \rangle$

$\langle \text{余留无符号数} \rangle \rightarrow d \langle \text{余留无符号数} \rangle \mid . \langle \text{十进小数} \rangle \mid e \langle \text{指数部分} \rangle \mid \varepsilon$

$\langle \text{十进小数} \rangle \rightarrow d \langle \text{余留十进小数} \rangle$

$\langle \text{余留十进小数} \rangle \rightarrow e \langle \text{指数部分} \rangle \mid d \langle \text{余留十进小数} \rangle \mid \varepsilon$

$\langle \text{指数部分} \rangle \rightarrow d \langle \text{余留整指数} \rangle \mid s \langle \text{整指数} \rangle$

$\langle \text{整指数} \rangle \rightarrow d \langle \text{余留整指数} \rangle$

$\langle \text{余留整指数} \rangle \rightarrow d \langle \text{余留整指数} \rangle \mid \varepsilon$

其中s表示正或负号。

如 25.55e+5 和 2.1

# 正规式

---

正规式也称正则表达式,正规表达式 (regular expression) 是说明单词的模式(pattern)的一种重要的表示法 (记号), 是定义正规集的数学工具。我们用以描述单词符号。

下面是正规式和它所表示的正规集的递归定义。

# 正规式的定义

---

定义（正规式和它所表示的正规集）：

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

1.  $\varepsilon$ 和 $\Phi$ 都是 $\Sigma$ 上的正规式，它们所表示的正规集分别为 $\{\varepsilon\}$ 和 $\{\}$ ；

# 正规式的定义

---

- 2。任何  $a \in \Sigma$ ,  $a$  是  $\Sigma$  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  $\Sigma$  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,  $(e_1)$ ,  $e_1 \mid e_2$ ,  $e_1 \bullet e_2$ ,  $e_1^*$  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  $\Sigma$  上的正规式, 仅由这些正规式所表示的集合才是  $\Sigma$  上的正规集。

# 正规式中的符号

---

其中的“|”读为“或”（也有使用“+”代替“|”的）；“•”读为“连接”；“\*”读为“闭包”（即，任意有限次的自重复连接）。在不致混淆时，括号可省去，但规定算符的优先顺序为“\*”、“•”、“|”。连接符“•”一般可省略不写。“\*”、“•”和“|”都是左结合的。



## 例子3.2

令  $\Sigma = \{a, b\}$ ,  $\Sigma$  上的正规式和相应的正规集的例子有:

正规式

正规集

$a$

$\{a\}$

$a \mid b$

$\{a, b\}$

$ab$

$\{ab\}$

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

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

$a^*$

$\{\epsilon, a, a, \dots \text{任意个 } a \text{ 的串}\}$

$(a \mid b)^*$

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

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

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

# 正规式举例

## 例3.3

令 $\Sigma = \{l, d\}$ , 则 $\Sigma$ 上的正规式  $r = l(l \mid d)^*$  定义的正闭集 $L(r)$ :  
 $\{l, ll, ld, ldd, \dots\}$ , 其中 $l$ 代表字母,  $d$ 代表数字, 正规式 即是  
字母(字母|数字)\*, 它表示的正规集中的每个元素的模式  
是“字母打头的字母数字串”, 就是Pascal和 多数程序设计  
语言允许的标识符的词法规则.

## 例3.2

$\Sigma = \{d, \bullet, e, +, -\}$ ,

则 $\Sigma$ 上的正规式  $d^*(\bullet dd^* \mid \varepsilon)(e(+ \mid - \mid \varepsilon)dd^* \mid \varepsilon)$  表示的是  
无符号数的集合。其中 $d$ 为0~9的数字。

**程序设计语言的单词都能用正规式 来定义.**

- 
- 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 $e_1=e_2$ 。

- 例如:  $e_1 = (a \mid b)$ ,  $e_2 = b \mid a$

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

• 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.  $\varepsilon r = r, r\varepsilon = r$   $\varepsilon$ 是“连接”的恒等元素  
零一律
- 6.  $r \mid r = r$   
 $r^* = \varepsilon \mid r \mid rr \mid \dots$  “或”的抽取律

# 正规文法和正规式的等价性

---

一个正规语言可以用正规文法定义，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。

对每个正规式，存在一个生成同一语言的正规文法。

# 正规式 $\rightarrow$ 正 $\square$ 文法

问题描述：将 $\Sigma$ 上的一个正规式 $r$ 转换成文法 $(V_N, V_T, P, S)$ 。

转换方法：令 $V_T = \Sigma$ ，确定产生式和 $V_N$ 的元素用如下方法：

1. 选择一个非终结符号 $S$ 生成类似的正规式产生式的形式 $S \rightarrow r$ 。
2. 若 $x$ 和 $y$ 都是正 $\square$ 式：
  - 若 $A \rightarrow xy$ 的正 $\square$ 式 $\square$ 生式，重写成 $A \rightarrow xB, B \rightarrow y$ ， $B$ 是新 $\square$ 的， $B \in V_N$ 。
  - 对形如 $A \rightarrow x^*y$ 的正 $\square$ 式 $\square$ 生式，重写 $\square A \rightarrow xA, A \rightarrow y$
  - 若形如 $A \rightarrow x | y$ 的正 $\square$ 式 $\square$ 生式，重写 $\square A \rightarrow x, A \rightarrow y$
3. 不断利用上述 $\square$ 做 $\square$ ，直到每个 $\square$ 生式都符合正 $\square$ 文法的格式

# 正规式 $\rightarrow$ 正规文法例

例3.4：将 $r=a(a|d)^*$ 转换成相应的正规文法。

令 $S$ 是文法的开始符号，

首先形成 $S \rightarrow a(a|d)^*$ ,

然后形成 $S \rightarrow aA, A \rightarrow (a|d)^*$ ,

再 $\square\square$ 成 $S \rightarrow aA, A \rightarrow (a|d)A, A \rightarrow \varepsilon$

进而变换成符合正规文法产生式的形式

$\rightarrow aA, A \rightarrow aA, A \rightarrow dA, A \rightarrow \varepsilon$

$S$

# 正规文法→正□式

基本上是正规式→正□文法的逆□程，最后只剩下一个开始符号定□的正□式，□□□□如下：

	文法产生式	正规式
规则1	$A \rightarrow xB, B \rightarrow y$ (必须包含A和B的所有产生式)	$A=xy$
规则2	$A \rightarrow xA \mid y$ (必须包含A的所有产生式)	$A=x^*y$
规则3	$A \rightarrow x, A \rightarrow y$	$A=x y$



# 正规文法→正则式例

例3.5：文法 $G[S]$ :  $S \rightarrow aA$ ,  $S \rightarrow a$ ,  $A \rightarrow aA$ ,  $A \rightarrow dA$ ,  $A \rightarrow a$ ,  
 $A \rightarrow d$  □ □ □ 正则式

首先有： $S = aA | a$       $A = (aA | dA) | (a | d)$

将A的正则式 □ □ □  $A = (a|d)A | (a|d)$ ，又 □ □ □：  
 $A = (a|d)^*(a|d)$ ，再将A的右端代入S的正则式得：  
 $S = a(a|d)^*(a|d) | a$

再利用正则式的代数 □ □ 可依次得到

$S = a((a|d)^*(a|d) | \epsilon)$       $S = a(a|d)^*$

即  $a(a|d)^*$  为所求。

## 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)。

# 关于有穷自动机我们将讨论如下题目

---

确定的有穷自动机DFA

不确定的有穷自动机NFA

NFA的确定化

DFA的最小化

# 确定的有穷自动机DFA

一个确定的有穷自动机（DFA） $M$ 是一个五元组： $M = (K, \Sigma, f, S, Z)$  其中

1.  $K$ 是一个有穷集，它的每个元素称为一个状态；
2.  $\Sigma$ 是一个有穷字母表，它的每个元素称为一个输入符号，所以也称 $\Sigma$ 为输入符号表；
3.  $f$ 是转换函数，是在 $K \times \Sigma \rightarrow K$ 上的映射，即，如  $f(k_i, a) = k_j$ ， $(k_i \in K, k_j \in K)$  就意味着，当前状态为 $k_i$ ，输入符为 $a$ 时，将转换为下一个状态 $k_j$ ，我们把 $k_j$ 称作 $k_i$ 的一个后继状态；
4.  $S \in K$ 是唯一的一个初态；
5.  $Z \subset K$ 是一个终态集，终态也称可接受状态或结束状态。

# 一个DFA 的例子：

---

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

$$f(S, a) = U$$

$$f(V, a) = U$$

$$f(S, b) = V$$

$$f(V, b) = Q$$

$$f(U, a) = Q$$

$$f(Q, a) = Q$$

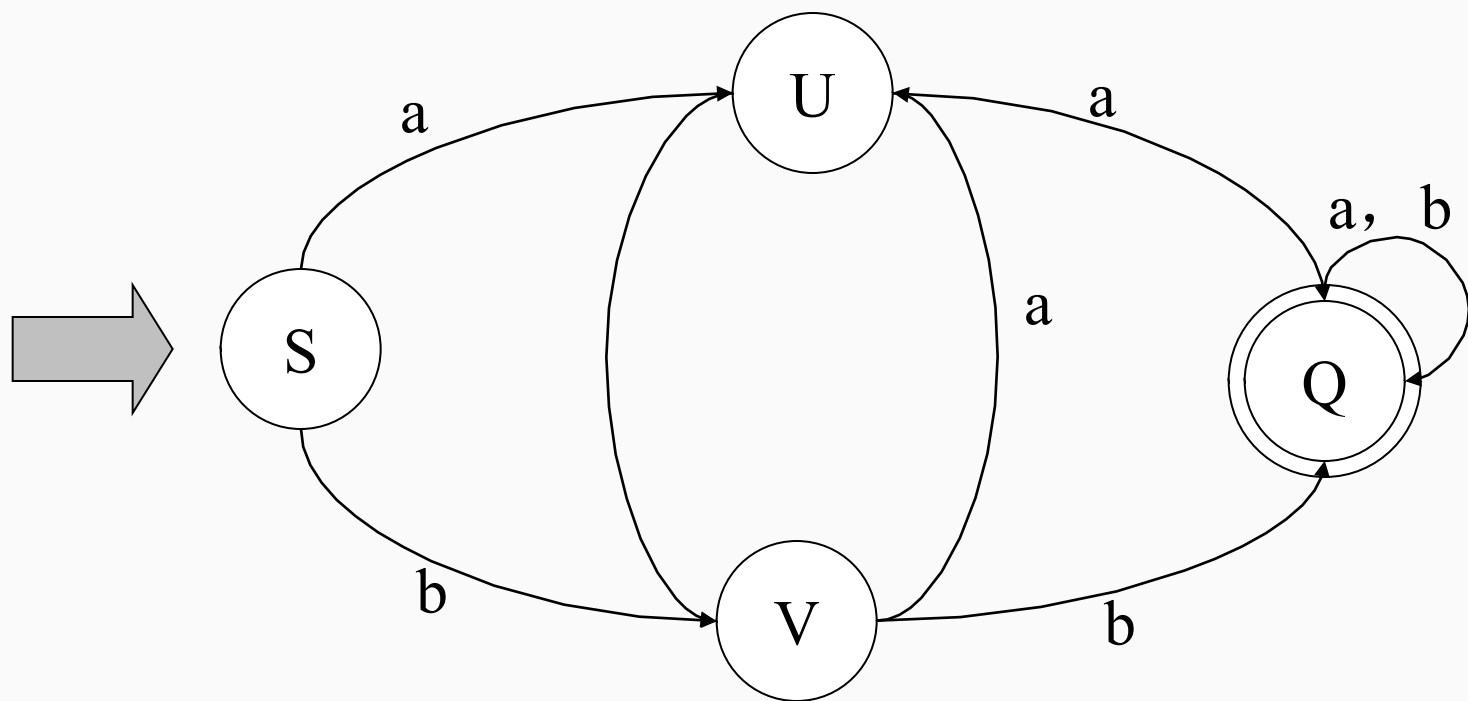
$$f(U, b) = V$$

$$f(Q, b) = Q$$

# 状态图

- 一个DFA可以表示成一个状态图(或称状态转换图)。
- 假定DFA  $M$  含有 $m$ 个状态， $n$ 个输入字符，那么这个状态图含有 $m$ 个结点，每个结点最多有 $n$ 个弧射出。
- 整个图含有唯一一个初态结点和若干个终态结点，初态结点冠以双箭头“ $\Rightarrow$ ”或标以“-”，终态结点用双圈表示或标以“+”；
- 若  $f(k_i, a) = k_j$ ，则从状态结点 $k_i$ 到状态结点 $k_j$ 画标记为 $a$ 的弧。

# DFA 的状态图表示



# 矩阵表示

---

- 一个DFA还可以用一个矩阵表示
- 该矩阵的行表示状态，列表示输入字符，矩阵元素表示相应状态行和输入字符列下的新状态，即 $k$ 行 $a$ 列为 $f(k,a)$ 的值。
- 用双箭头“ $\Rightarrow$ ”标明初态；否则第一行即是初态，相应终态行在表的右端标以1，非终态标以0。



# DFA的矩阵表示

状态 \ 字符	a	b	
S	U	V	0
U	Q	V	0
V	U	Q	0
Q	Q	Q	1

# 运行

## $\Sigma^*$ 上的符号串 $t$ 在DFA $M$ 上运行

一个输入符号串 $t$ , (将它表示成 $t_1t_x$ 的形式, 其中  $t_1 \in \Sigma$ ,  $t_x \in \Sigma^*$ ) 在DFA  $M = (K, \Sigma, f, S, Z)$  上运行的定义为:

$$f(Q, t_1t_x) = f(f(Q, t_1), t_x) \quad \text{其中 } Q \in K$$

扩充转换函数 $f$ 为  $K \times \Sigma^* \rightarrow K$  上的映射, 且:  $f(k_i, \varepsilon) = k_i$

# 识别（接受）

---

$\Sigma^*$ 上的符号串 $t$ 被DFA  $M$ 接受

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

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**识别（接受）**。

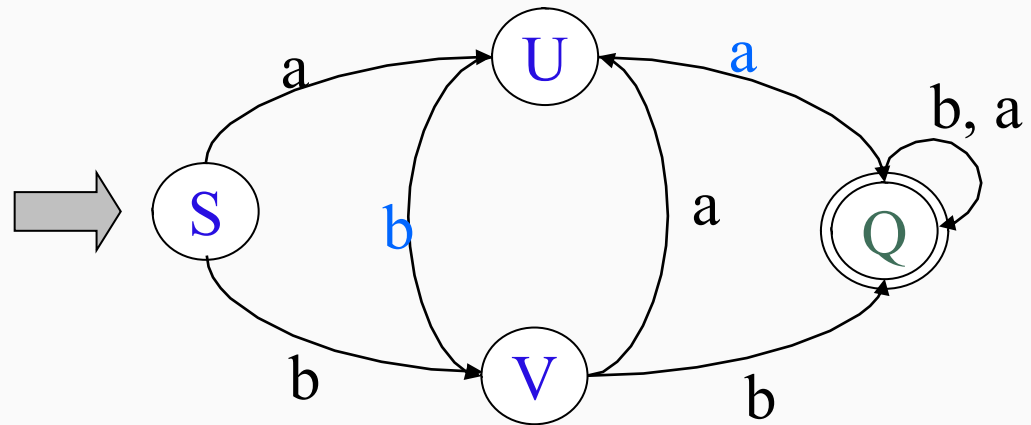
即, 若存在一条从初态结点到某一终态结点的道路,  
且这条道路上所有弧的标记符连接成的符号串等于 $t$ ,  
则称 $t$ 为DFA  $M$ 所**识别（接受）**。

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。



---

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:

$\Sigma$ 上一个符号串集 $V \subset \Sigma^*$ 是正闭的, 当且仅当存在一个 $\Sigma$ 上的确定有穷自动机 $M$ , 使得

$V=L(M)$ 。

# DFA的确定性

---

- 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 $\Sigma$ 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。

# DFA的行为的程序模拟.

---

**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')**

# 不确定的有穷自动机NFA

---

## 定义

NFA  $M = \{K, \Sigma, f, S, Z\}$ , 其中  $K$  为状态的有穷非空集,  $\Sigma$  为有穷输入字母表,  $f$  为  $K \times \Sigma^*$  到  $K$  的子集 ( $2^K$ ) 的一种映射,  $S \subset K$  是一个非空的初始状态集,  $Z \subset K$  为终止状态集.



# NFA举例

---

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

其中

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

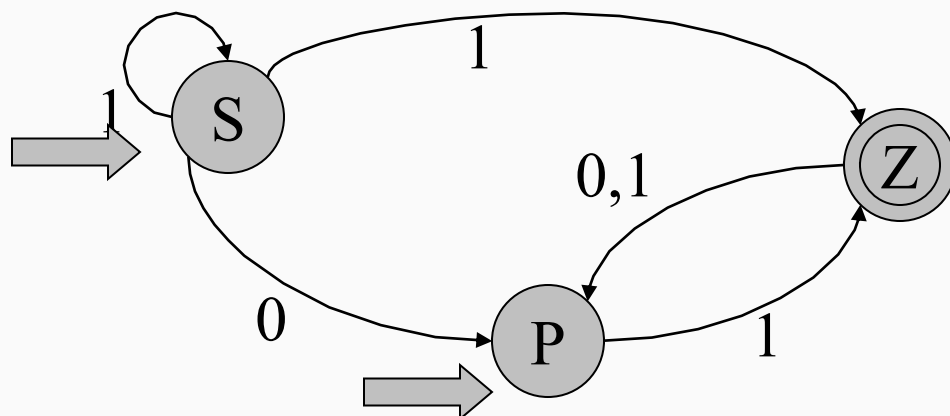
$f(Z, 0) = \{P\}$

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

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

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

## 状态图表示

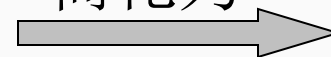


# 矩阵表示

## 矩阵表示

	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

## 类似DFA, 对NFA $M = \{K, \Sigma, f, S, Z\}$ 也有如下定义

- $\Sigma^*$ 上的符号串 $t$ 在NFA  $M$ 上运行..  
一个输入符号串 $t$ , (我们将它表示成 $t_1t_x$ 的形式, 其中 $t_1 \in \Sigma$ ,  $t_x \in \Sigma^*$ ) 在NFA  $M$ 上**运行**的定义为:  
 $f(Q, t_1t_x) = f(f(Q, t_1), t_x)$  其中 $Q \in K$ .
- $\Sigma^*$ 上的符号串 $t$ 被NFA  $M$ 接受  
若 $t \in \Sigma^*$ ,  $f(S_0, t) = P$ , 其中 $S_0 \in S$ ,  $P \in Z$ ,  
则称 $t$ 为NFA  $M$ 所**接受** (识别)

$\Sigma^*$ 上的符号串 $t$ 被NFA  $M$ 接受也可以这样理解

- 对于 $\Sigma^*$ 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(不理睬那些标记为 $\varepsilon$ 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 $\varepsilon$ ，那么空字可为 $M$ 所接受。

000

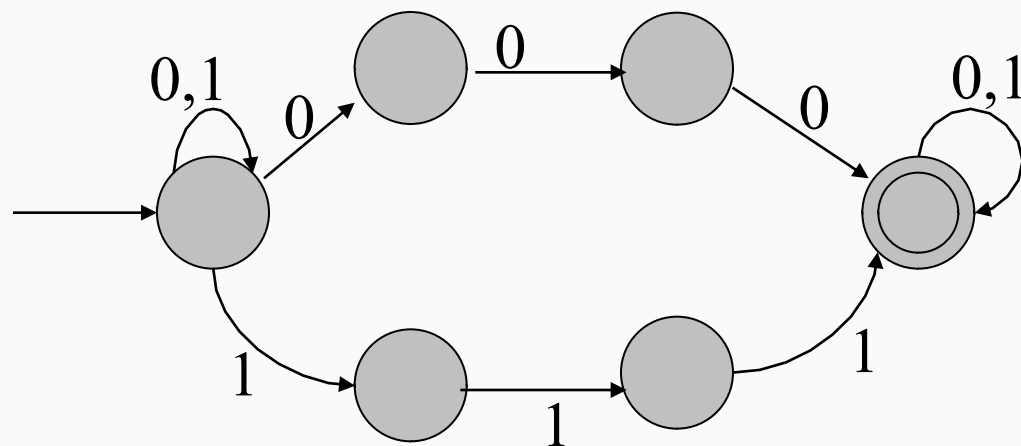
111

1010001

110000001

00

01100



---

NFA  $M$ 所能接受的符号串的全体记为

$L(M)$

结论：

$\Sigma$ 上一个符号串集  $V \subset \Sigma^*$  是正闭的，当且仅当存在一个  $\Sigma$  上的不确定的有穷自动机  $M$ ，使得  $V = L(M)$ 。

# DFA和NFA的关系

---

DFA是NFA的特例. 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  $M$ 。

有一种算法，将NFA转换成接受同样语言的DFA. 这种算法称为子集法.

与某一NFA等价的DFA不唯一.



# NFA确定化基本思路

---

从NFA的矩阵表示中可以看出，表项通常是一状态的集合，而在DFA的矩阵表示中，表项是一个状态，NFA到相应的DFA的构造的基本思路是：**DFA的每一个状态对应NFA的一组状态。**  
DFA使用它的一个状态去记录在NFA读入一个输入符号后可能达到的所有状态。

# 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 $\Sigma$ ；
  - 3 转换函数是这样定义的：  
$$d([S_1 S_2 \dots S_j], a) = [R_1 R_2 \dots R_t] \quad \text{其中} \quad \{R_1, R_2, \dots, R_t\} = \varepsilon\text{-closure}(\text{move}(\{S_1, S_2, \dots, S_j\}, a))$$
  - 4  $S_0 = \varepsilon\text{-closure}(K_0)$ 为M的开始状态；
  - 5  $S_t = \{[S_i S_k \dots S_e], \text{ 其中 } [S_i S_k \dots S_e] \in S \text{ 且 } \{S_i, S_k, \dots, S_e\} \cap K_t \neq \Phi\}$

# 定义对状态集合I的几个有关运算

---

1. 状态集合I的  $\epsilon$ -闭包，表示为  $\epsilon\text{-closure}(I)$ ，定义为一状态集，是状态集I中的任何状态S经任意条  $\epsilon$  弧而能到达的状态的集合。  
状态集合I的任何状态S都属于  $\epsilon\text{-closure}(I)$ 。
2. 状态集合I的  $a$  弧转换，表示为  $\text{move}(I, a)$  定义为状态集合J，其中J是所有那些可从I中的某一状态经过一条  $a$  弧而到达的状态的全体。

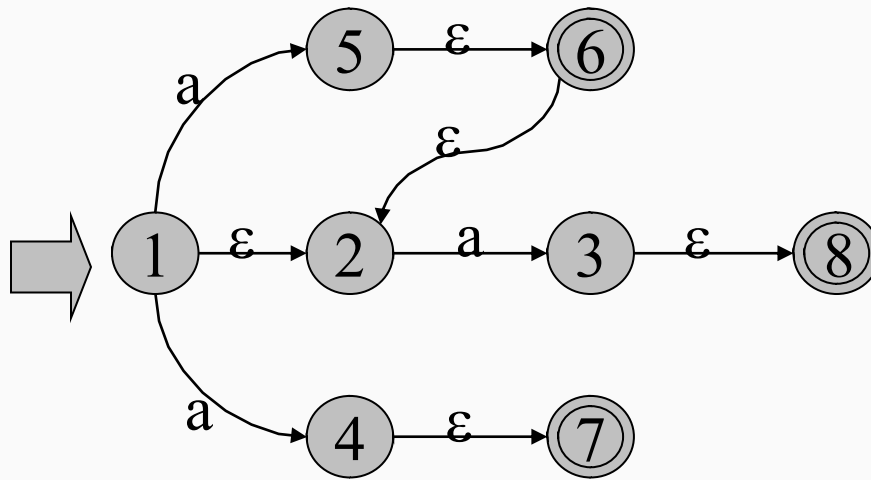
## 状态集合I的有关运算的例子

$I = \{1\}$ ,  $\varepsilon\text{-closure}(I) = \{1, 2\}$ ;

$I = \{5\}$ ,  $\varepsilon\text{-closure}(I) = \{5, 6, 2\}$ ;

$\text{move}(\{1, 2\}, a) = \{5, 3, 4\}$

$\varepsilon\text{-closure}(\{5, 3, 4\}) = \{2, 3, 4, 5, 6, 7, 8\}$ ;



## 构造NFA $N$ 的状态 $K$ 的子集的算法:

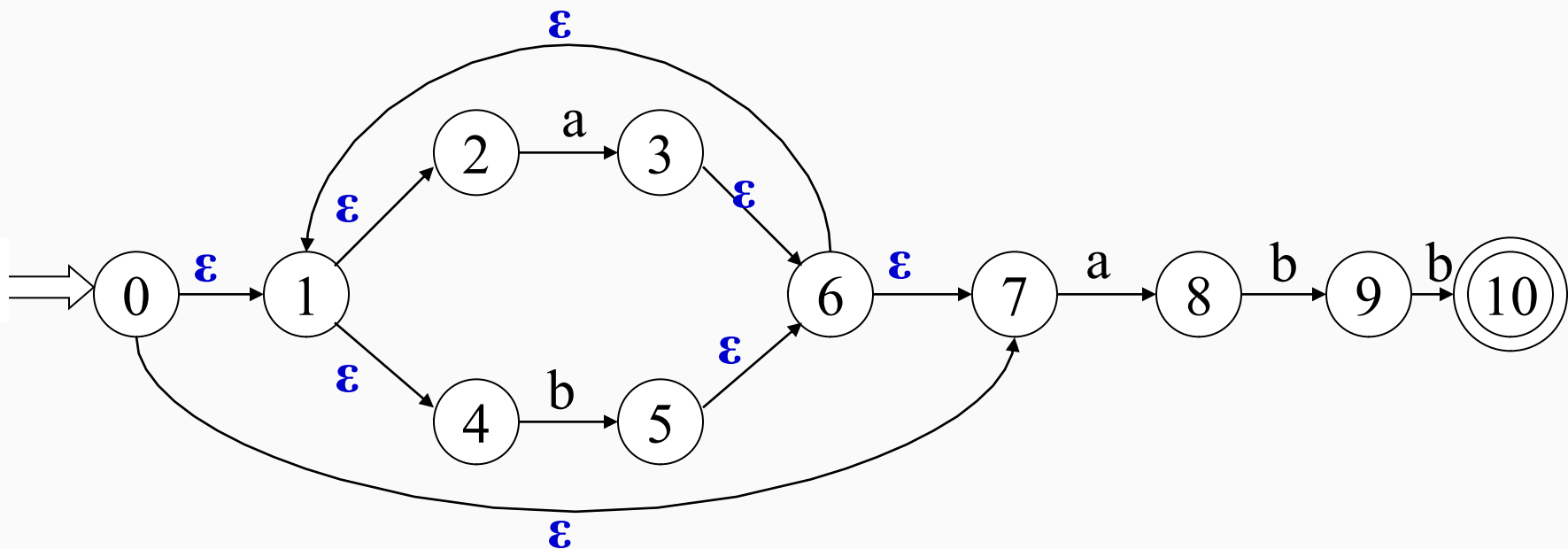
假定所构造的子集族为 $C$ , 即 $C = (T_1, T_2, \dots, T_I)$ , 其中 $T_1, T_2, \dots, T_I$ 为状态 $K$ 的子集。

- 1 开始, 令 $\varepsilon\text{-closure}(K_0)$ 为 $C$ 中唯一成员, 并且它是未被标记的。

---

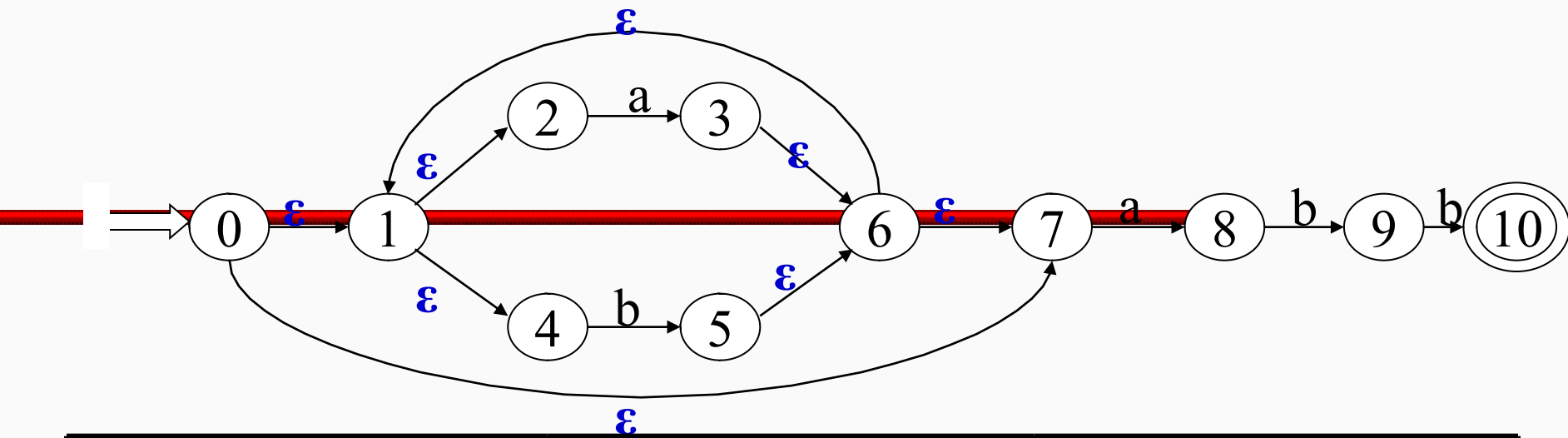
```
2  while (C中存在尚未被标记的子集T) do
    {
        标记T;
        for 每个输入字母a do
            {
                U:=  $\epsilon$ -closure(move(T,a));
                if U不在C中 then
                    将U作为未标记的子集加在C中
            }
    }
```

# NFA确定化举例



**NFA N**

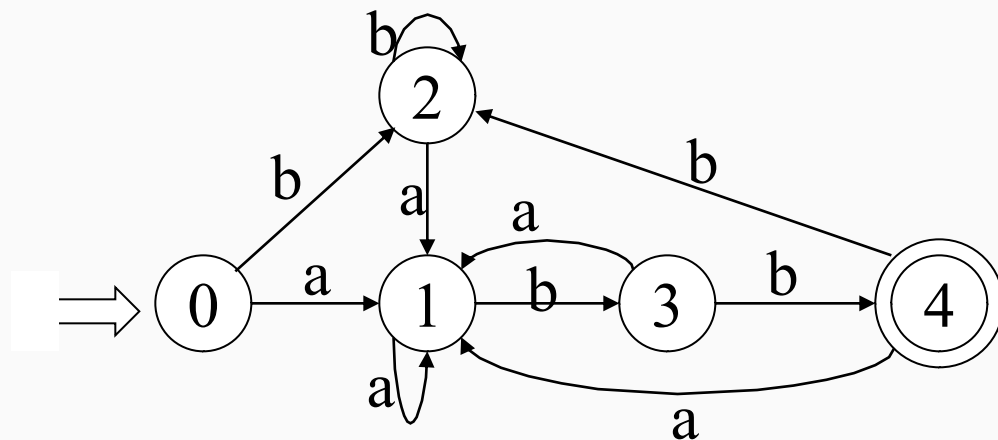




	$\epsilon$ -closure(move( $T_i, a$ ))	$\epsilon$ -closure(move( $T_i, b$ ))	
$T_0 = \epsilon$ -closure(0) $= \{0, 1, 2, 4, 7\}$ <b>=&gt;0</b>	$\{1, 2, 3, 4, 6, 7, 8\}$ 加入为 $T_1$ <b>1</b>	$\{1, 2, 4, 5, 6, 7\}$ 加入为 $T_2$ <b>2</b>	0
$T_1 = \{1, 2, 3, 4, 6, 7, 8\}$ <b>1</b>	$\{1, 2, 3, 4, 6, 7, 8\}$ 已存在 $T_1$ <b>1</b>	$\{1, 2, 4, 5, 6, 7, 9\}$ 加入为 $T_3$ <b>3</b>	0
$T_2 = \{1, 2, 4, 5, 6, 7\}$ <b>2</b>	$\{1, 2, 3, 4, 6, 7, 8\}$ 已存在 $T_1$ <b>1</b>	$\{1, 2, 4, 5, 6, 7\}$ 已存在 $T_2$ <b>2</b>	0
$T_3 = \{1, 2, 4, 5, 6, 7, 9\}$ <b>3</b>	$\{1, 2, 3, 4, 6, 7, 8\}$ 已存在 $T_1$ <b>1</b>	$\{1, 2, 4, 5, 7, 10\}$ 加入为 $T_4$ <b>4</b>	0
$T_4 = \{1, 2, 4, 5, 7, 10\}$ <b>4</b>	$\{1, 2, 3, 4, 6, 7, 8\}$ 已存在 $T_1$ <b>1</b>	$\{1, 2, 4, 5, 6, 7\}$ 已存在 $T_2$ <b>2</b>	1

# 最终生成DFA

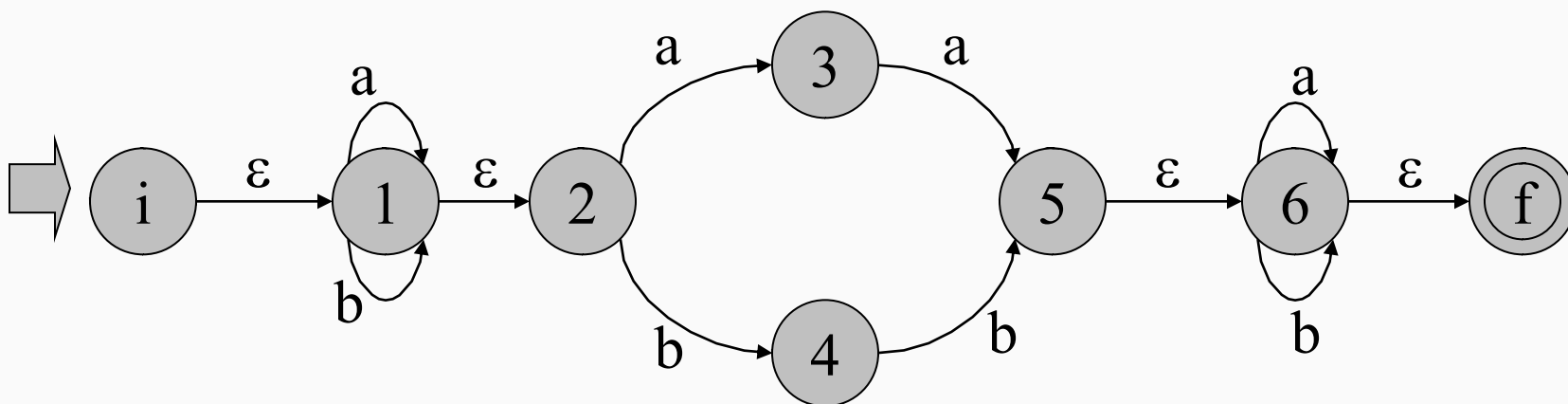
状态	a	b	
0	1	2	0
1	1	3	0
2	1	2	0
3	1	4	0
4	1	2	1

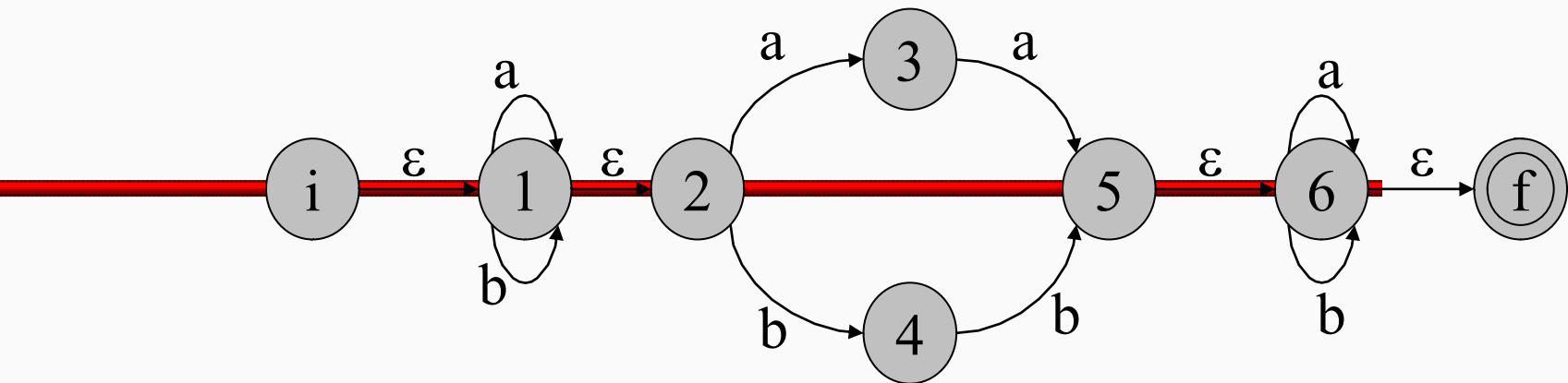


**DFA M**

# NFA的确定化

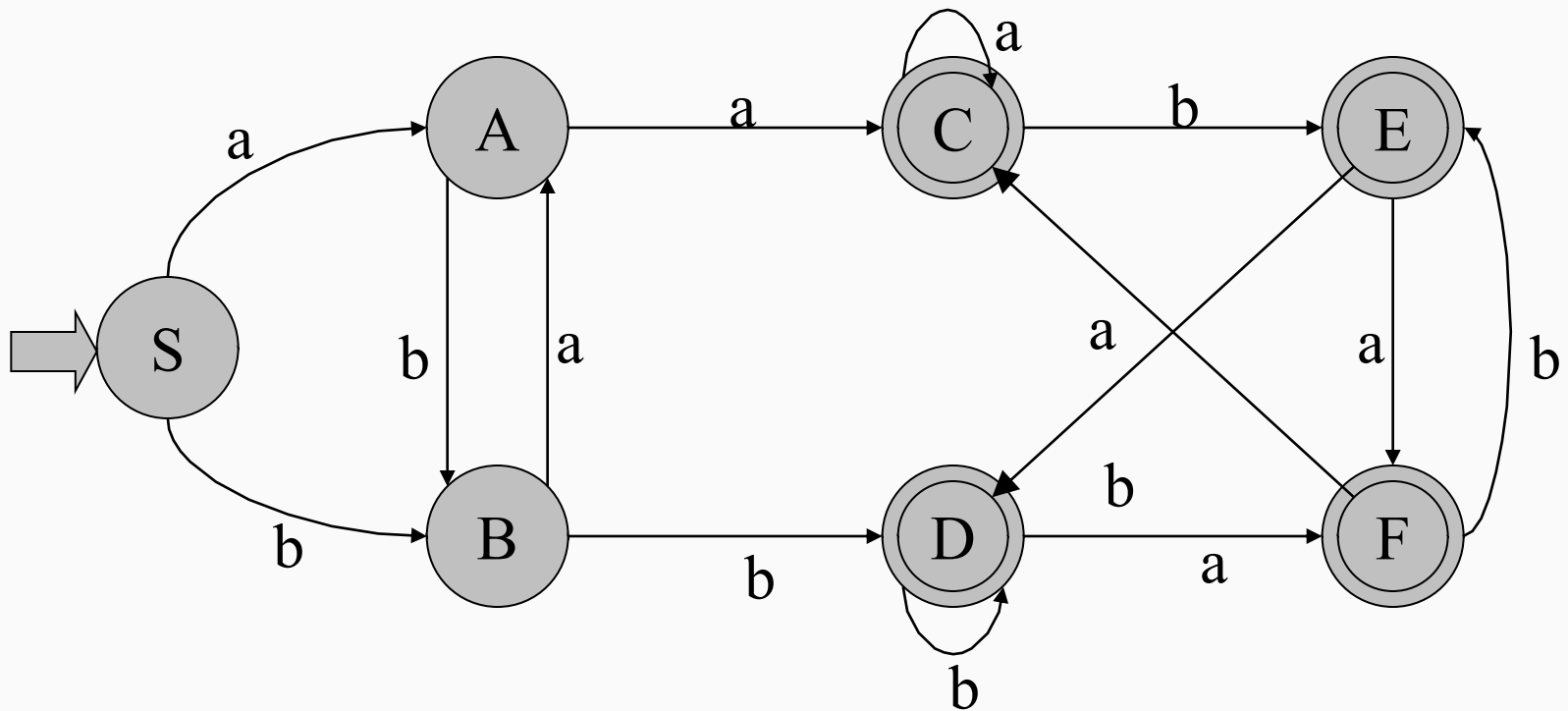
## 例子





0  
0  
0  
1  
1  
1  
1

# 等价的DFA



# 确定有穷自动机的化简

- 说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，是指这样的状态：从自动机的开始状态出发，任何输入串也不能到达的那个状态；或者从这个状态没有通路到达终态。

# 消除多余状态举例

	0	1		0	1		0	1			
s0	s1	s5	0	s0	s1	s5	0	s0	s1	s5	0
s1	s2	s7	1	s1	s2	s7	1	s1	s2	s7	1
s2	s2	s5	1	s2	s2	s5	1	s2	s2	s5	1
s3	s5	s7	0	s3	s5	s7	0	s3	s5	s7	0
s4	s5	s6	0	s4	s5	s6	0	s5	s3	s1	0
s5	s3	s1	0	s5	s3	s1	0	s7	s0	s1	1
s6	s8	s0	1	s6	s8	s0	1				
s7	s0	s1	1	s7	s0	s1	1				
s8	s3	s6	0	s8	s3	s6	0				

# DFA的最小化就是寻求最小状态DFA

---

最小状态DFA的含义:

没有多余状态(死状态)

没有两个状态是互相等价（不可区别）

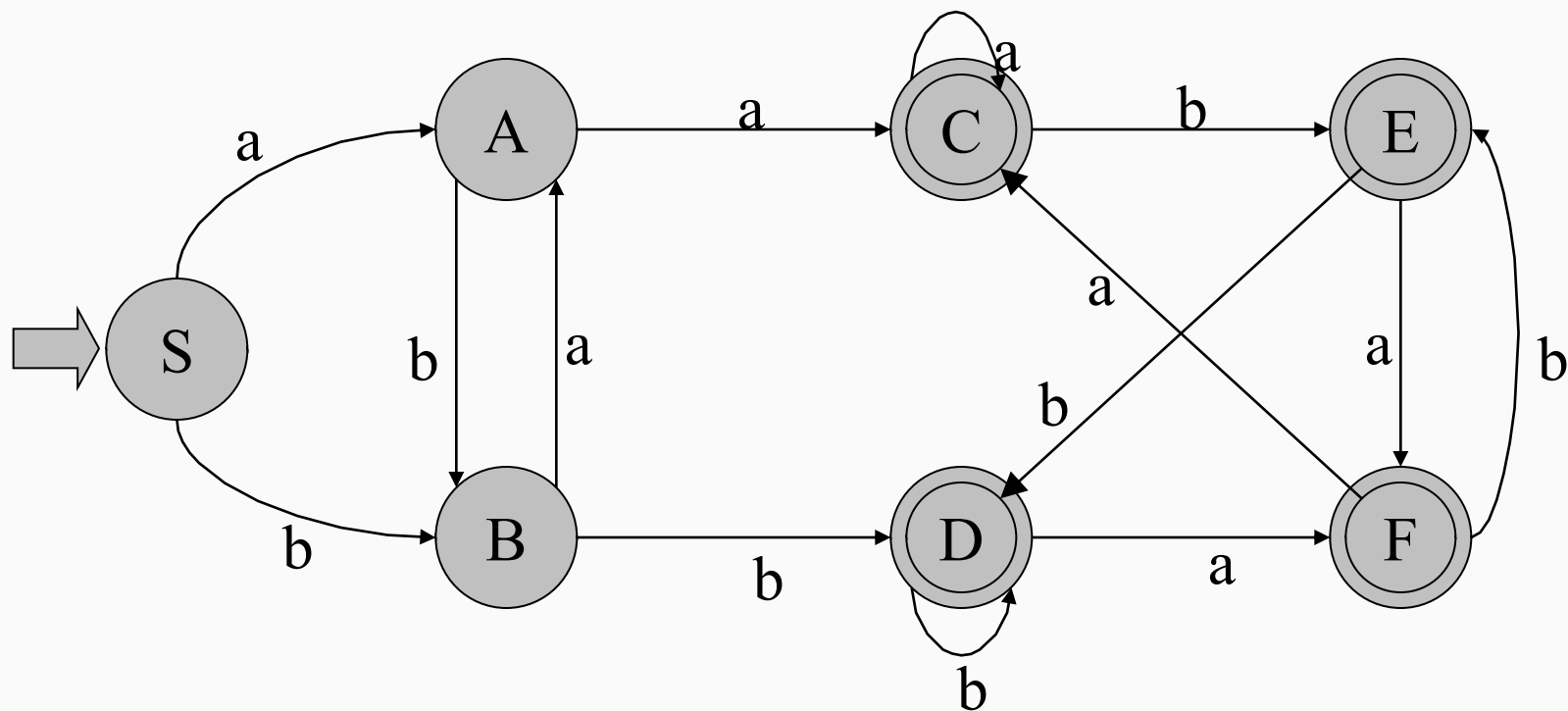
两个状态s和t等价的条件是:

一致性条件——同时为终态或者非终态

蔓延性条件——对于所有输入符号，状态s和状态t必须转换到等价的状态里。



# 等价状态举例



C、F等价，D、E等价，C、D等价

# 最小状态DFA

---

对于一个DFA  $M = (K, \Sigma, f, k_0, k_t)$ , 存在一个最小状态DFA  $M' = (K', \Sigma, f', k'_0, k'_t)$ , 使  $L(M') = L(M)$ .

## 结论

- 接受L的最小状态有穷自动机不计同构是唯一的。

# “分割法”

---

- DFA的最小化算法的核心

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。

# DFA的最小化算法

DFA  $M = (K, \Sigma, f, k_0, k_t)$ , 最小状态DFA  $M'$

1. 构造状态的一初始划分 $\Pi$ :

终态 $k_t$  和非终态 $K - k_t$ 两组(group) ;

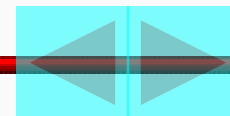
2. 对 $\Pi$ 施用过程PP构造新划分 $\Pi_{\text{new}}$

PP过程是这样的:

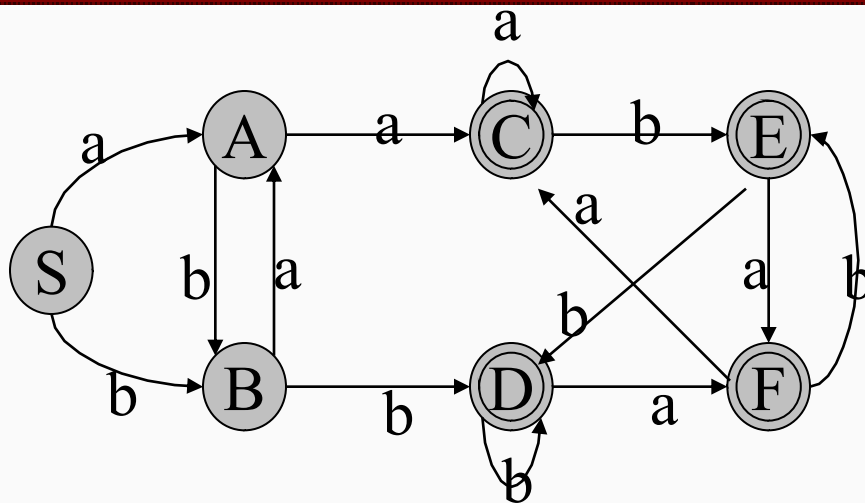
① 将 $\Pi$ 中的每一组 $G$ 划分为子组, 而 $G$ 中的两个状态 $k_i$ 和 $k_j$ 在同一子组中的充分必要条件是: 对所有输入符号 $a$ 而言,  $k_i$ 和 $k_j$ 转换后的状态都处于 $\Pi$ 中同样的组里;

② 形成的子组构成了新划分 $\Pi_{\text{new}}$ 中的 $G$ 。

3. 如  $\Pi_{\text{new}} = \Pi$ , 则令  $\Pi_{\text{final}} = \Pi$  并继续步骤4, 否则  $\Pi := \Pi_{\text{new}}$  重复2。
4. 为  $\Pi_{\text{final}}$  中的每一组选一代表, 这些代表构成  $M'$  的状态。若  $k$  是一代表且  $f(k, a) = t$ , 令  $r$  是  $t$  组的代表, 则  $M'$  中有一转换  $f'(k, a) = r$ 。  $M'$  的开始状态是含有  $S_0$  的那组的代表,  $M'$  的终态是含有  $F$  的那组的代表;
5. 去掉  $M'$  中的死状态。



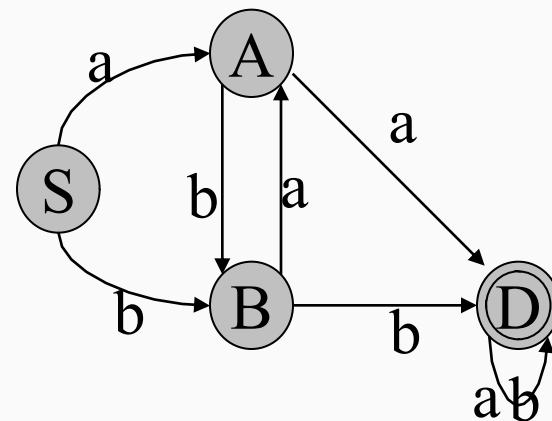
# DFA的最小化—例子



$\Pi_0: \{S, A, B\} \quad \{C, D, E, F\}$

$\Pi_1: \{A\} \quad \{S, B\} \quad \{C, D, E, F\}$

$\Pi_2: \{A\} \quad \{S\} \quad \{B\} \quad \{C, D, E, F\}$



## 3.5 正规式和有穷自动机的等价性

---

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

# 有穷自动机→正规式

## 问题描述：

如何为 $\Sigma$ 上的NFA  $M$ 构造相应的正规式 $r$

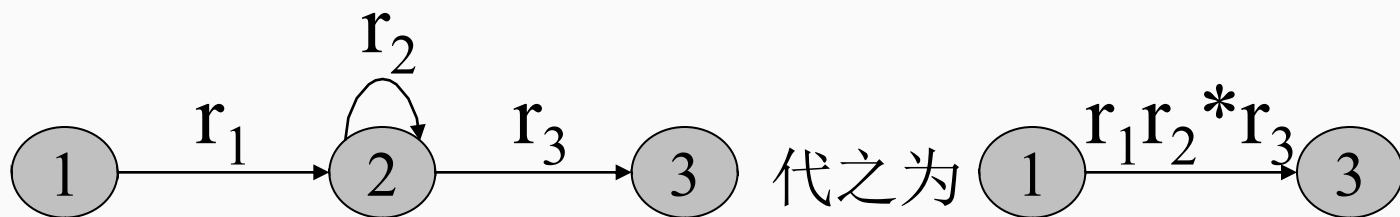
## 算法：

把状态转换图的概念拓广，令每条弧可用一个正规式作标记。

1. 在 $M$ 的状态图上加进两个结点，一个为 $x$ 结点，一个为 $y$ 结点。从 $x$ 结点用 $\varepsilon$ 弧连接到 $M$ 的所有初态结点，从 $M$ 的所有终态结点用 $\varepsilon$ 弧连接到 $y$ 结点。形成一个与 $M$ 等价的 $M'$ ， $M'$ 只有一个初态 $x$ 和一个终态 $y$ 。



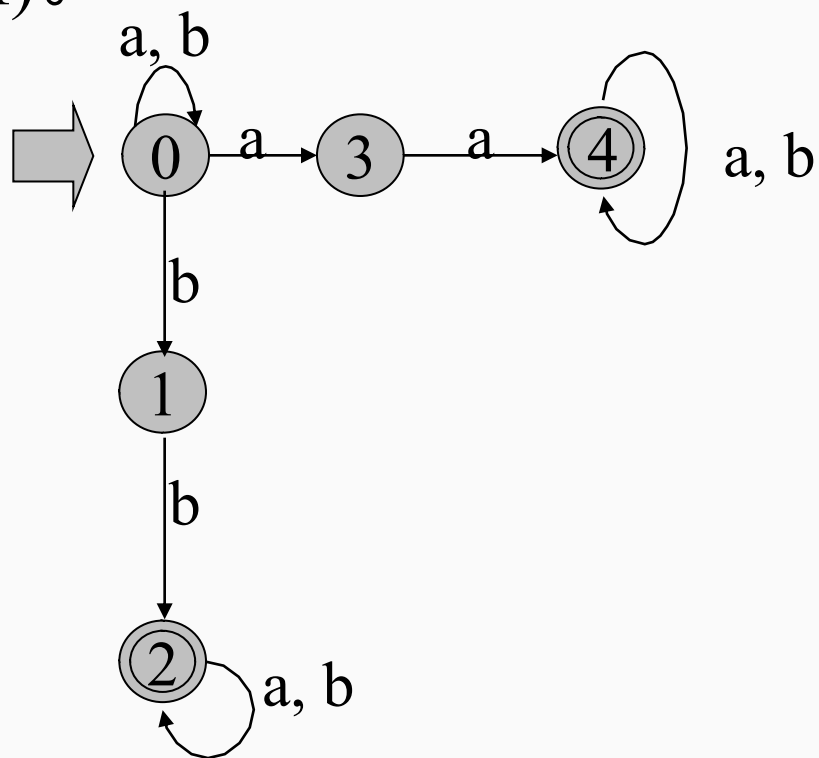
2. 逐步消去M'中的所有结点，直至只剩下x和y结点。在消结过程中，逐步用正规式来标记弧。其消结的规则如下：



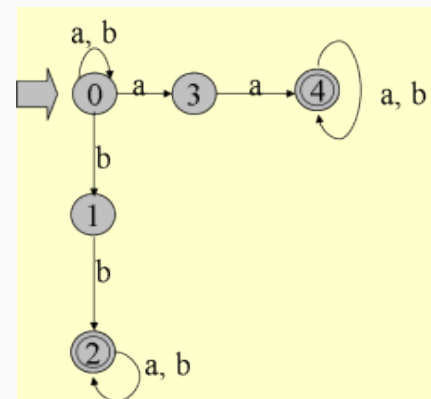
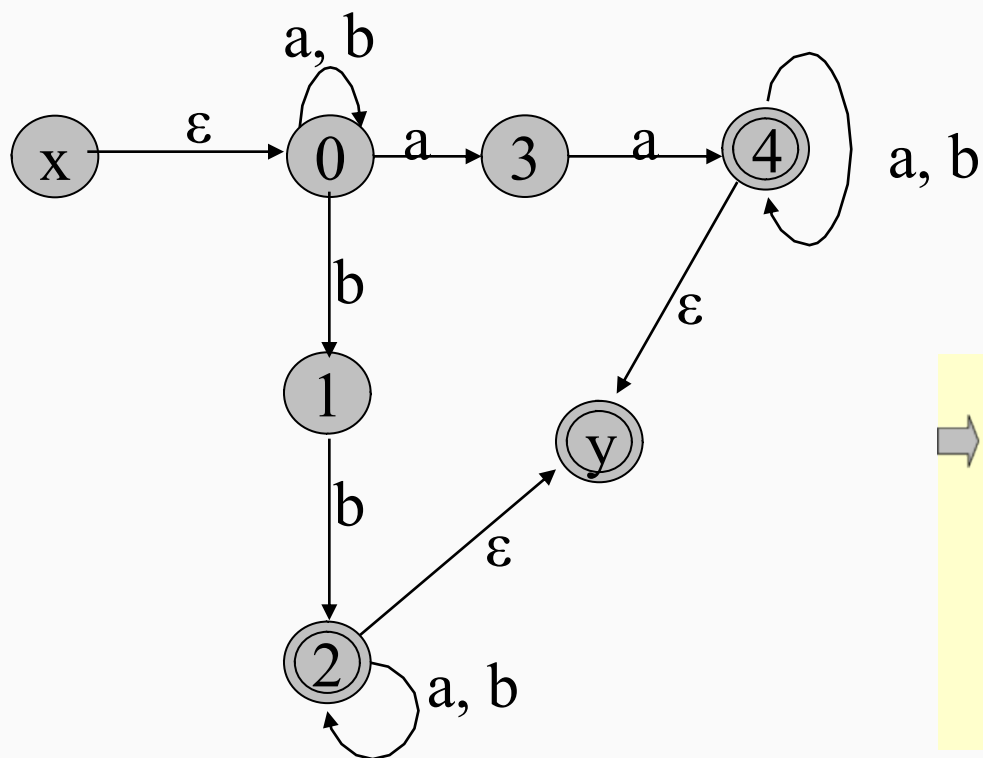
---

3. 最后 $x$ 和 $y$ 结点间的弧上的标记则为所求的正规式 $r$ 。

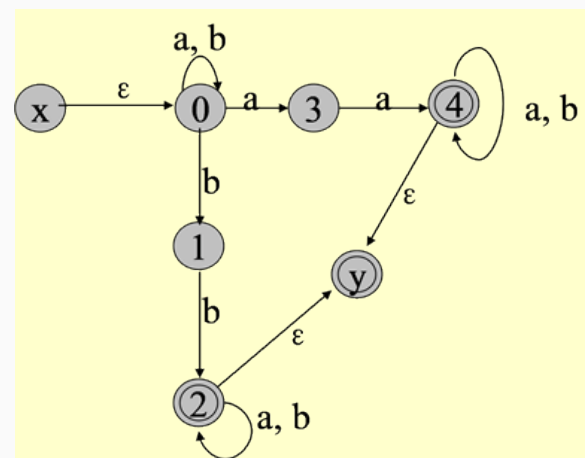
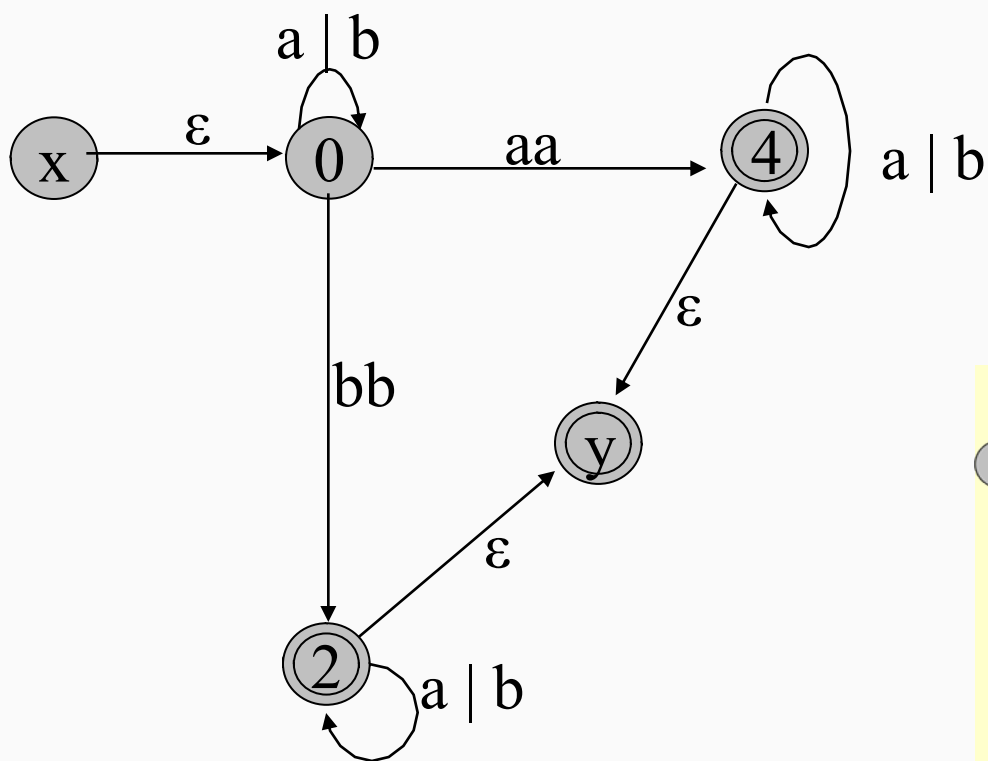
例：将下图中的NFA  $M$ ，求正规式 $r$ ，使得 $L(r)=L(M)$ 。



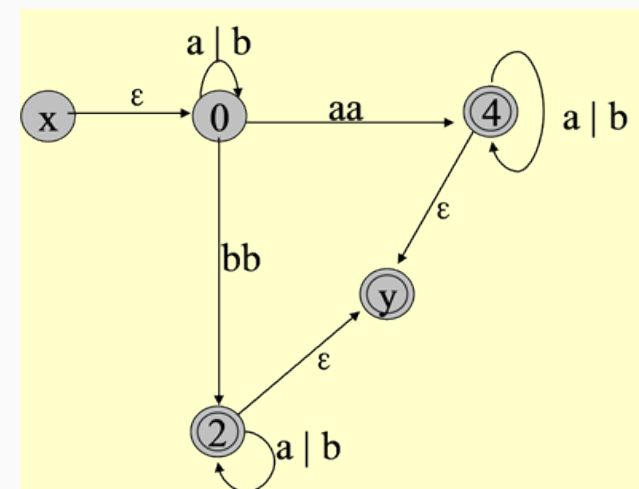
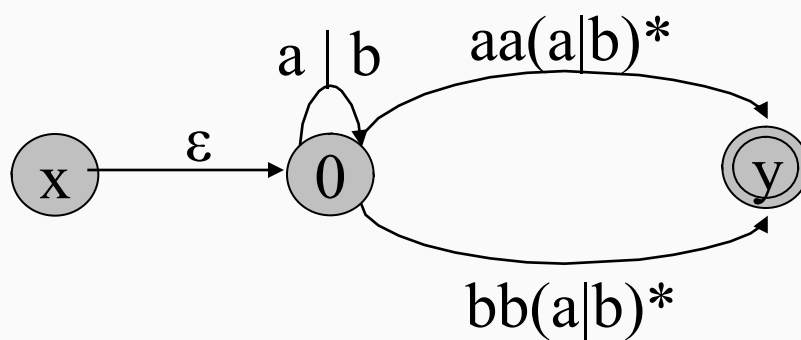
1. 加x和y结点，形成如下的自动机M'



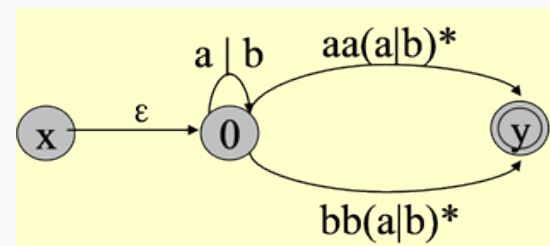
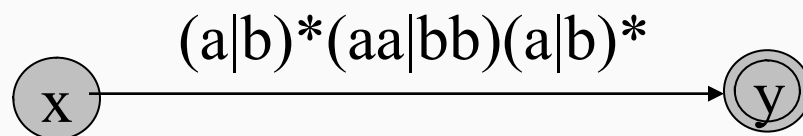
## 2. 消去M'中的1和3结点



### 3. 消去2和4结点



4. 消去0结点，只剩下x和y结点如下图：



$r=(a|b)^*(aa|bb)(a|b)^*$ 即为所求。

# 正␣式→有穷自动机

---

- 从 $\Sigma$ 上的一个正规式 $R$ 构造 $\Sigma$ 上的一个NFA  $M$ , 使得 $L(M)=L(R)$ 的方法。
- “语法制导”的方法, 即按正规式的语法结构指引构造过程, 构造规则具体描述如下:



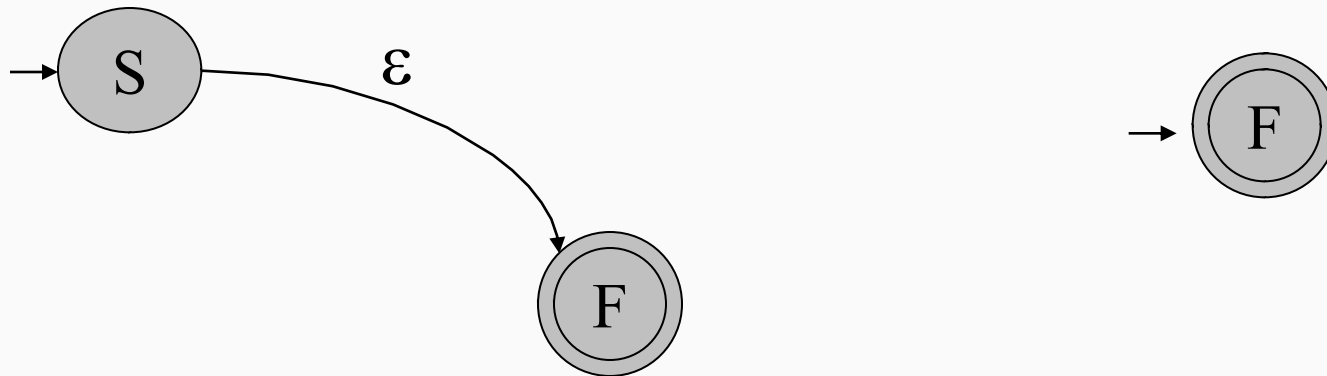
# 对于正规式 $R=\emptyset$ ,构造的NFA

---



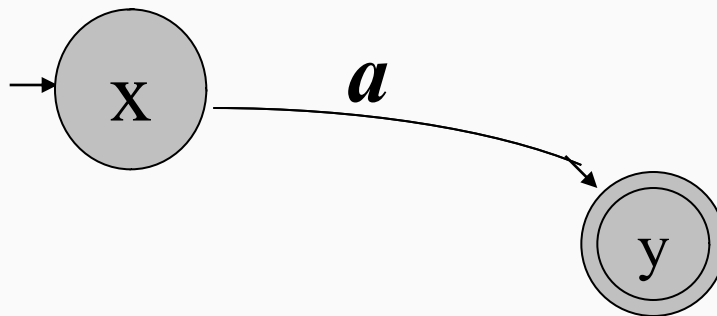
# 对于正规式 $\varepsilon$ ,构造的NFA (二种)

---



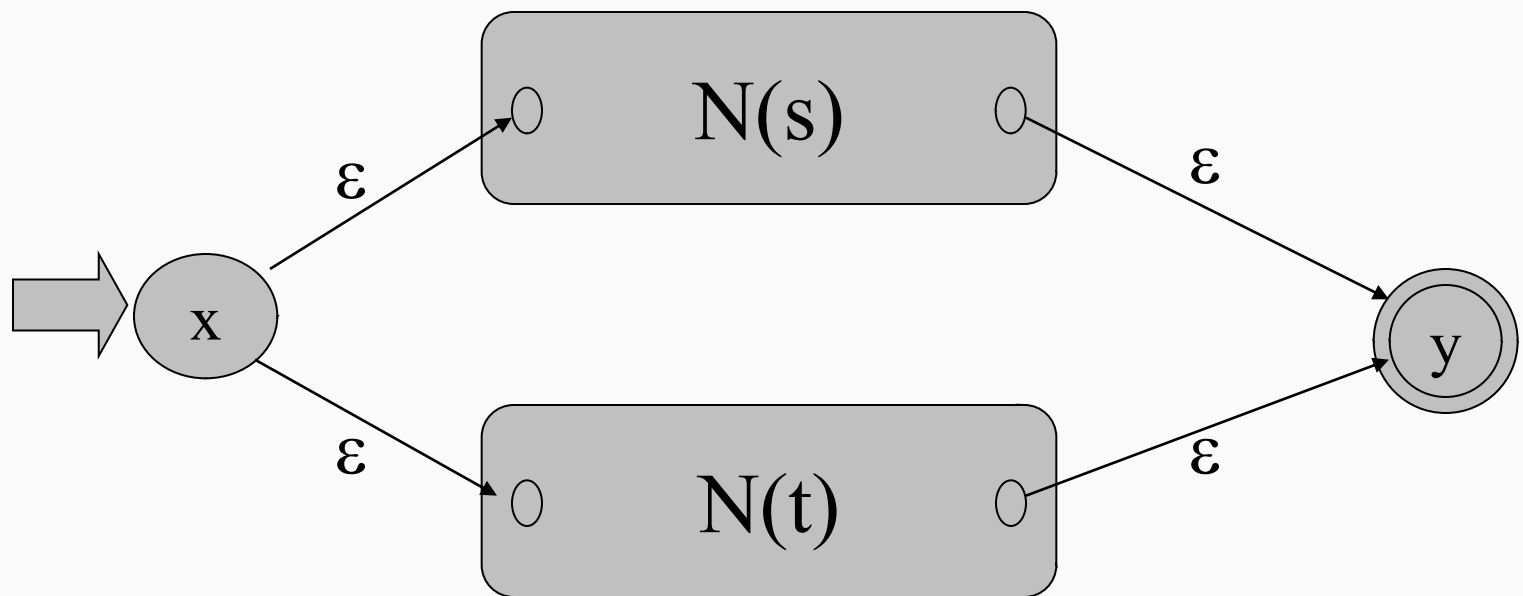
# 对于正规式 $a$ , $a \in \Sigma$ 构造的NFA

---



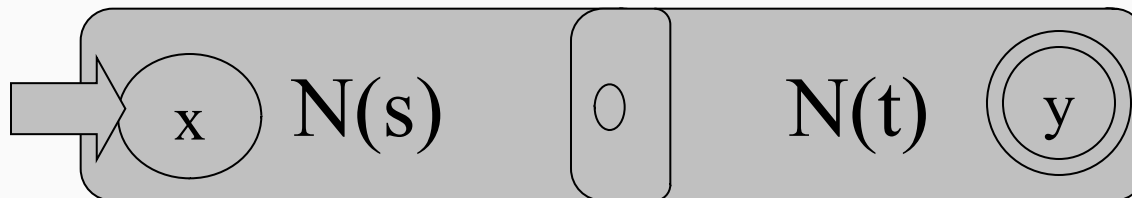
对于正规式 $r$ ,  $r = s|t$ 构造的NFA

---

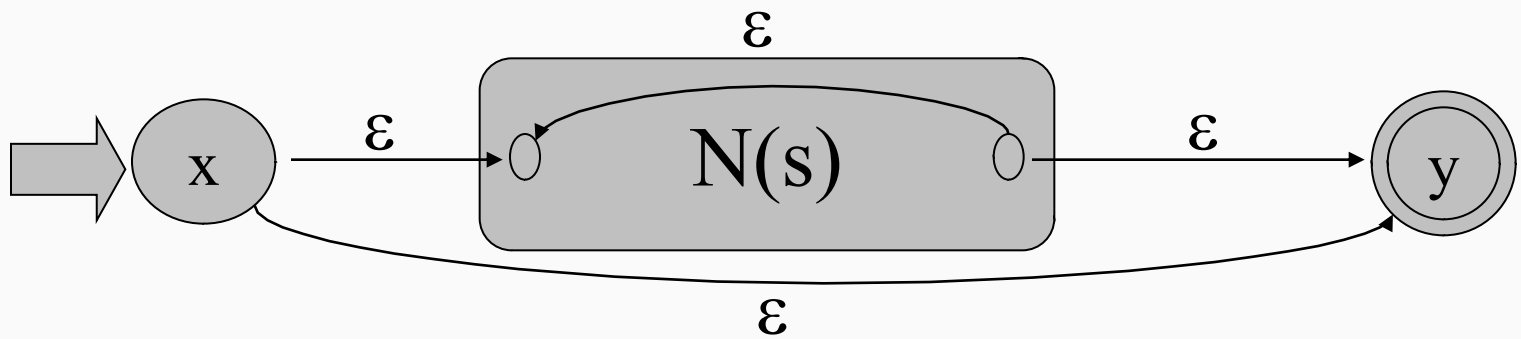


# 对于正规式 $r$ , $r=st$ 构造的NFA

---

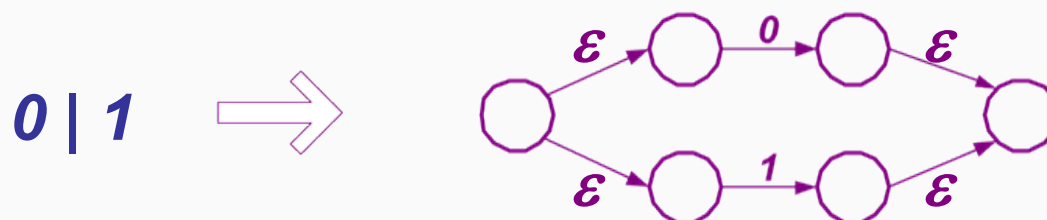
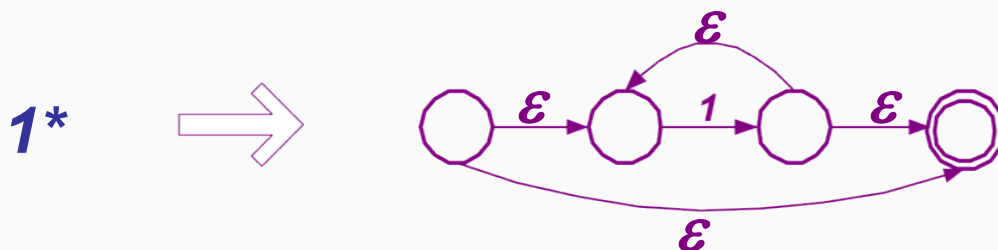


# 对于正规式 $r$ , $r=s^*$ 构造的NFA



# 举例:从正规表达式构造等价的 $\epsilon$ - NFA

正规表达式  $1^*0(0 \mid 1)^*$

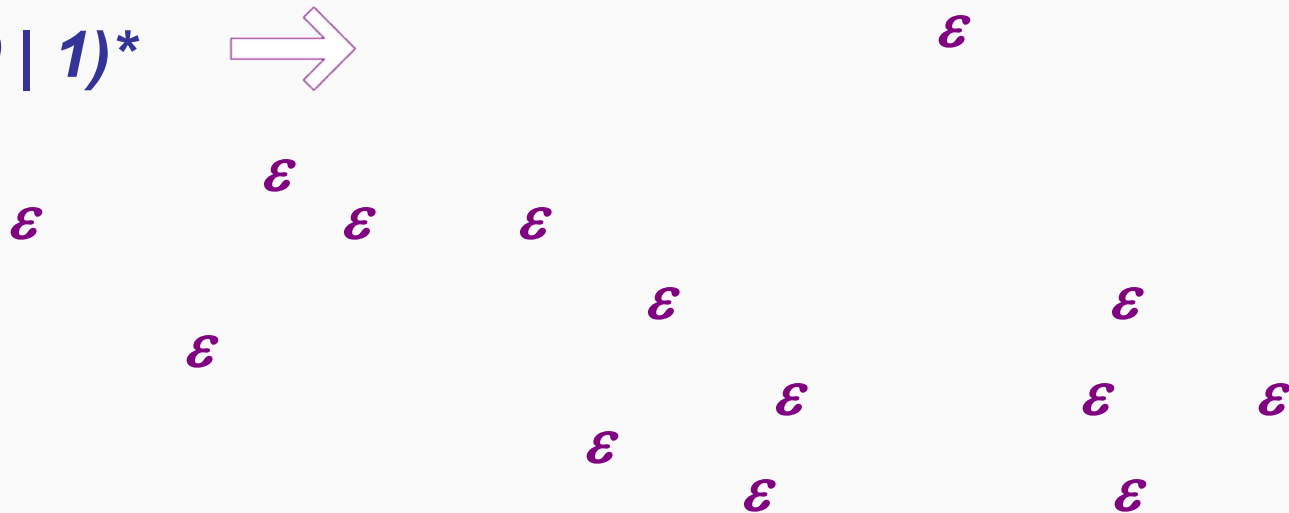


# 从正规表达式构造等价的 $\varepsilon$ - NFA

$(0 \mid 1)^*$



$1^*0(0 \mid 1)^*$

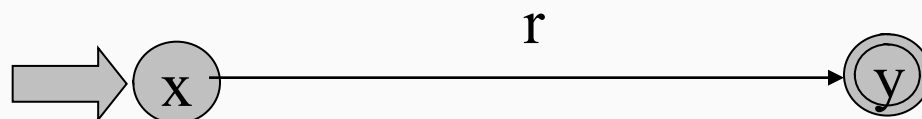




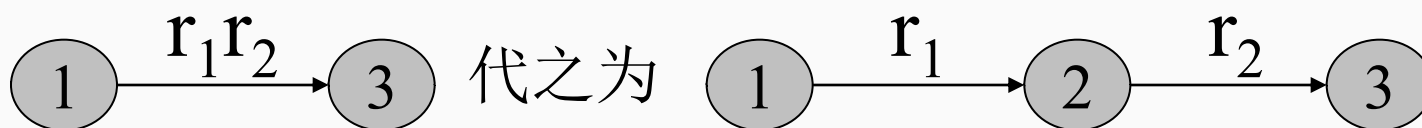
# 补充：正规式 $\rightarrow$ 有穷自动机的简易方法

---

1. 对于正规式 $r$ ，生成有穷自动机 $M$ 如下：



2. 逐步消去M中的所有标记的长度大于1的弧，  
规则如下：

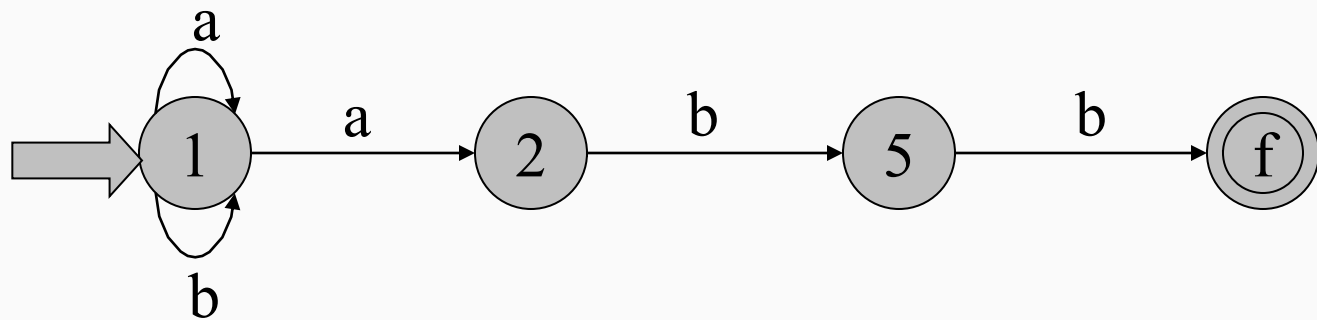


---

3. 所有弧上的标记都为 $a(a \in V_N)$ 或者 $\varepsilon$ 。所得到的NFA即为所求。

例：  $R=(a|b)^* abb$

---



## 3.6 正规文法和有穷自动机的等价性

采用下面的规则可以从正规文法 $G$ 直接构造一个有穷自动机NFA  $M$ ，使得 $L(M)=L(G)$ 。

1.  $M$ 的字母表与 $G$ 的终结集相同；
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 $A \rightarrow tB$ 的□□(其中 $t$ □□□符号或 $\varepsilon$ ， $A$ 和 $B$ 为非终结符的产生式)，构造 $M$ 的一个□□函数 $f(A, t)=B$ ；
5. □ $G$ 中的形如 $A \rightarrow t$ 的□生式，构造 $M$ 的一个□□函数 $f(A, t)=Z$ 。

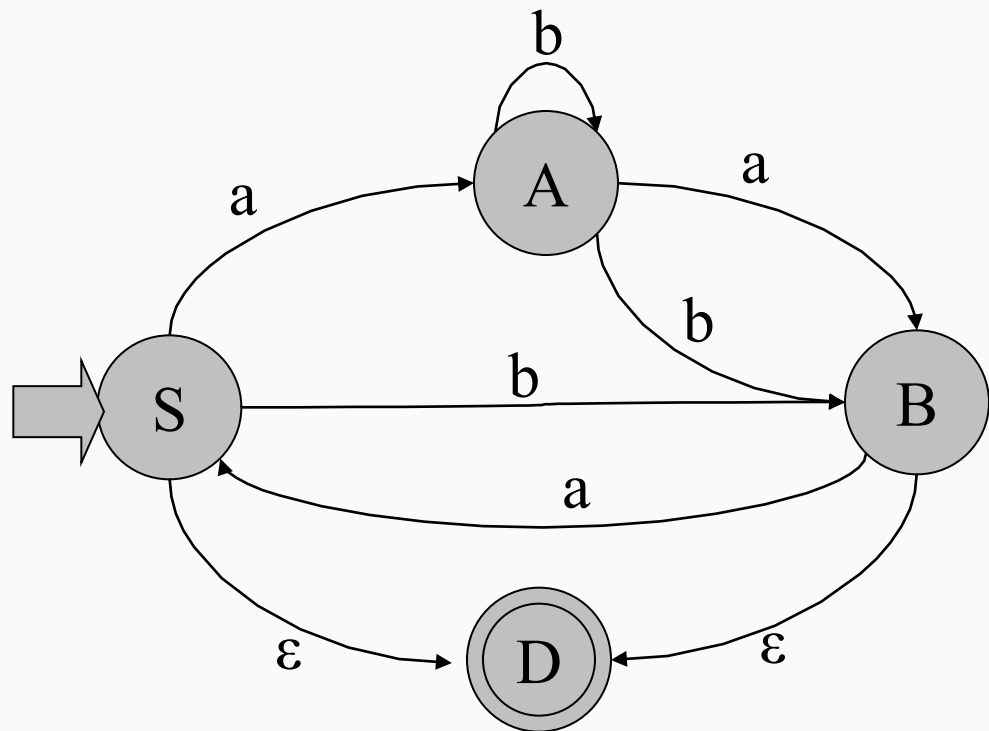
# 例：

构造与下述文法  
 $G[S]$ 等价的NFA  $M$ 。

$G[S]: S \rightarrow aA \mid bB \mid \varepsilon$

$A \rightarrow aB \mid bA$

$B \rightarrow aS \mid bA \mid \varepsilon$



# 有穷自动机→正□文法

---

对转换函数 $f(A, t)=B$ 可写成一个产生式:

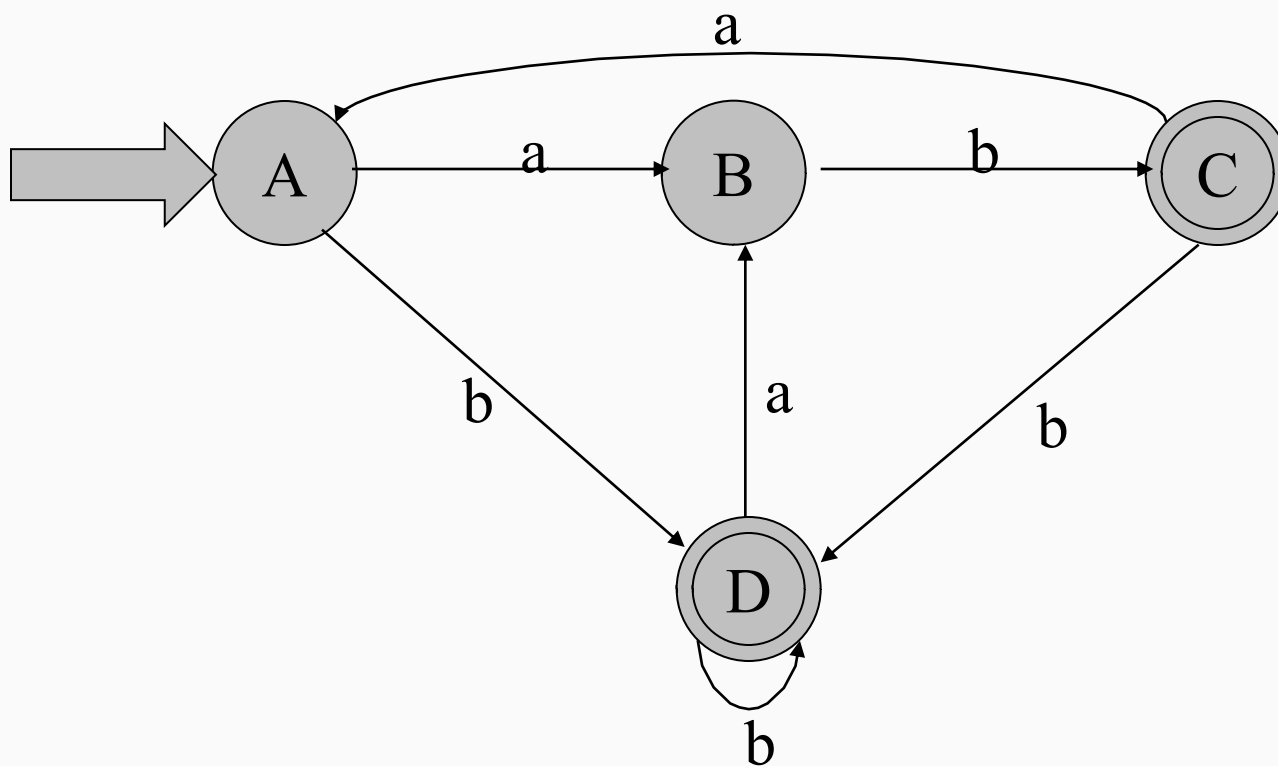
$$A \rightarrow tB$$

对可接受状态 $Z$ , 增加一产生式:

$$Z \rightarrow \varepsilon$$

# 例：

给出与下图的NFA等价的正规文法G





等价文法  $G = (\{A, B, C, D\}, \{a, b\}, P, A)$ , 其中  $P$  为:

$A \rightarrow aB$

$C \rightarrow \varepsilon$

$A \rightarrow bD$

$D \rightarrow aB$

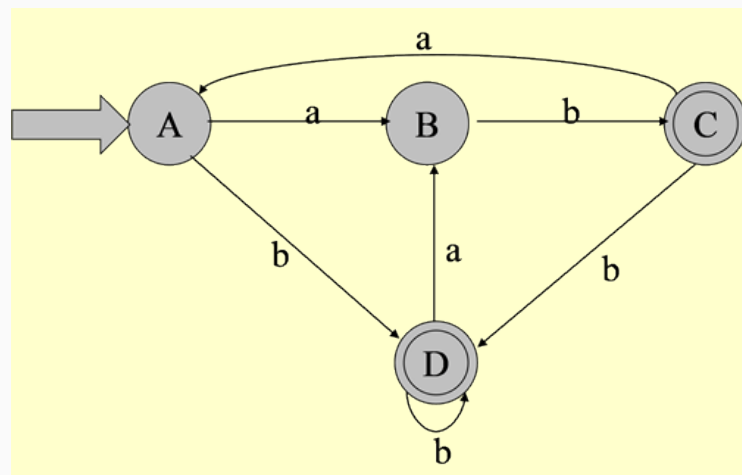
$B \rightarrow bC$

$D \rightarrow bD$

$C \rightarrow aA$

$D \rightarrow \varepsilon$

$C \rightarrow bD$



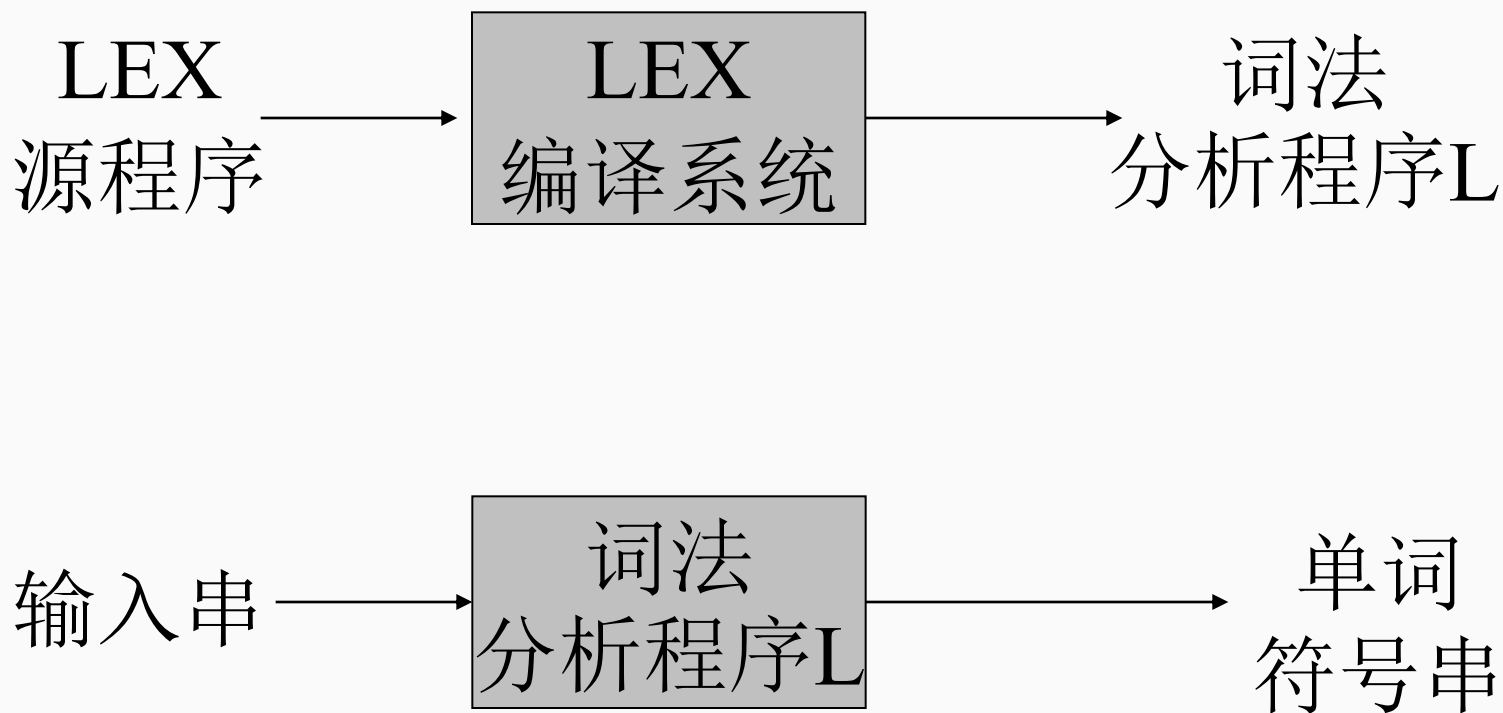
## 3.7 词法分析程序的自动构造工具

---

- 正规式用于说明(描述)单词的结构十分简洁方便。而把一个正规式编译(或称转换)为一个NFA进而转换为相应的DFA，这个NFA或DFA正是识别该正规式所表示的语言的句子的识别器。
- 基于这种方法来构造词法分析程序的工具很多
- LEX是一个广泛使用的工具。

# LEX编译系统的作用

---



# LEX程序组成

---

- 说明部分：包括变量的说明，常量说明和正规定义。
- 转换规则：包括正规式和识别出单词后词法分析器所应采取的动作。
- 辅助过程：这些过程可以分别编译并置于词法分析器中。

- 
- 词法分析程序的设计技术可应用于其它领域，比如查询语言以及信息检索系统等，这种应用领域的程序设计特点是，通过字符串模式的匹配来引发动作，回想LEX，说明词法分析程序的语言，可以看成是一个模式动作语言。
  - 词法分析程序的自动构造工具也广泛应用于许多方面，如用以生成一个程序，可识别印刷电路板中的缺陷，又如开关线路设计和文本编辑的自动生成等。

# 习题

---

1 (2), (3)

3

4

5

8

# 本章小结

---

词法分析程序是编译第一阶段的工作，它读入字符流的源程序，按照词法规则识别单词，交由语法分析程序接下去。

本章讲述了词法分析程序设计原则，并介绍了分别作为正规集的描述机制和识别机制的正规式和有穷动机。在此基础上给出了词法分析程序自动构造工具如LEX的原理。