

第四章 词法分析

词法分析程序亦称为扫描器

扫描器的任务是识别基本的语法单位
——单词

扫描器的输出是语法分析程序的输入

第四章 词法分析

词法分析程序的设计和实现

- 首先需要描述和刻画语言中的原子单位——单词，其次需要识别单词和执行某些相关的动作。描述程序设计语言的词法的机制是**3型文法**和正则表达式，识别机制是有穷状态自动机。



第四章 词法分析

- 设计词法分析程序
- 单词的描述工具
- 单词的识别系统

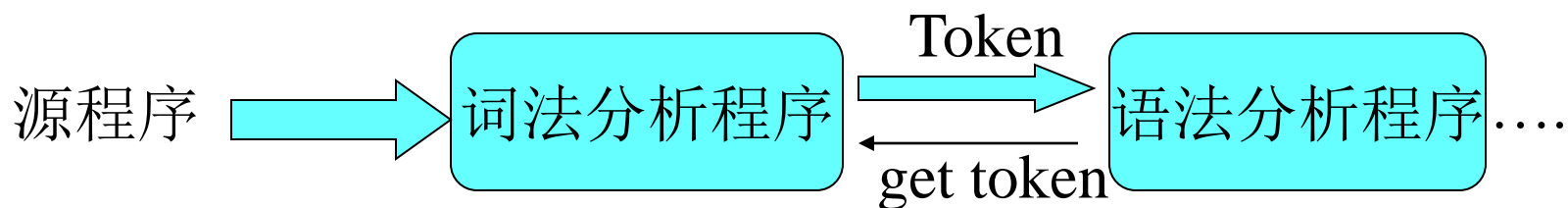
- 4.1 词法分析程序
- 4.2 正规表达式与正规集（正规语言）
- 4.3 有穷自动机
- 4.4 有穷自动机和正规表达式
- 4.5 有穷自动机和正规文法

4.1词法分析程序

■ 词法分析（lexical analysis）

- 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。

4.1词法分析程序



- 主要任务：
 - 读源程序，产生单词符号，并转换为token表示
- 其他任务：
 - 滤掉空格，删除注释、换行符
 - 对行列计数
 - 发现并定位词法错误，并尽量改正
 - 建立符号表、常数表等表格，

4.1词法分析程序

- 4.1.1词法分析器的输入缓冲区
- 4.1.2词法分析器的输出
- 单词符号一般可分为下列五种：
 - 基本字，关键字
 - 标识符
 - 常数（量）
 - 运算符
 - 界符

4.1词法分析程序

- 输出的二元式（token）表示：
- （单词种别，单词自身的值）
- Token的种类：
- 1、有些单词，只需要值，如基本字；
- 2、有些单词，还需要其他信息，如标示符。
- 例： $A:=B+2$
(Id的整数码,指向A的符号表的入口指针)
(运算符的整数码,':=')
(Id的整数码,指向B的符号表的入口指针)
(运算符的整数码,'+')
(常数的整数码, 2)

■ 4.1.3 以状态转换图为例设计词法分析器

4.2 正规表达式与正规集（正规语言）

- 正规表达式（**regular expression**）是描述单词符号的一种方便工具，也是定义正规集的工具。
- 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 $\Sigma' = \{\Phi, \varepsilon, |, \bullet, *, (,)\}$ 。
 - **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 $\{\}$ ；

4.2 正规表达式与正规集（正规语言）

- **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) , $e_1 \mid e_2$, $e_1 \bullet e_2$, 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 Σ 上的正规式, 仅由这些正规式所表示的字集才是 Σ 上的正规集。

4.2正规表达式与正规集（正规语言）

- 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，
- 规定算符的优先顺序为：

“（”、
高

“）”、
“*”、
“●”、
“|”

“)”、
低
- 连接符“●”一般可省略不写。
- “*”、“●”和“|”都是左结合的。

4.2 正规表达式与正规集（正规语言）

■ 例1 令 $\Sigma=\{a, b\}$, Σ 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^*	$\{\varepsilon, a, aa, \dots\}$ 任意个a的串

4.2 正规表达式与正规集（正规语言）

– 正规式

正规集

– $(a \mid b)^*$

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

– $(a \mid b)^*(aa \mid bb)(a \mid b)^*$ $\{\Sigma^*$ 上所有含有两个相继的 **a** 或两个相继的 **b** 组成的串}

4.2 正规表达式与正规集（正规语言）

- 例2 $\Sigma = \{l, d\}$, $r = l(l \mid d)^*$ 定义的正规集?
- $\{l, ll, ld, ldd, \dots\}$
- 例3 $\Sigma = \{d, ., e, +, -\}$, 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。

4.2 正规表达式与正规集（正规语言）

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

4.2 正规表达式与正规集（正规语言）

■ 设 r, s, t 为正规式，正规式服从的代数规律有：

- 1、 $r | s = s | r$ “或”服从交换律
- 2、 $r | (s | t) = (r | s) | t$ “或”的可结合律
- 3、 $(rs)t = r(st)$ “连接”的可结合律
- 4、 $r(s | t) = rs | rt$
 $(s | t)r = sr | tr$ 分配律

4.2 正规表达式与正规集（正规语言）

– 5、 $\varepsilon r = r, r\varepsilon = r$ ε 是“连接”的恒等元素
零一律

– 6、 $r \mid r = r$
 $r^* = \varepsilon \mid r^+$
 $r^+ = rr^*$

正规文法和正规式

■ 文法的定义

– $G = \{V_N, V_T, P, S\}$

– V_N : 非终结符的非空有穷集

– V_T : 终结符的非空有穷集

– P : 产生式的非空有穷集

$\alpha \rightarrow \beta \quad \alpha \in (V_N \cup V_T)^*$ 且至少含一个非
终结符, $\beta \in (V_N \cup V_T)^*$

– S : $\in V_N$, 称为开始符号

■ 正规文法: G 的任何产生式为 $A \rightarrow aB$ 或 $A \rightarrow a$,
其中 $a \in V_T \cup \varepsilon$, $A, B \in V_N$

正规文法和正规式

– $G1 = (\{S, A, B\}, \{a, b\}, S, P)$

其中 P : $S \rightarrow A$

$S \rightarrow B$

$A \rightarrow aA$

$A \rightarrow a$

$B \rightarrow bB$

$B \rightarrow b$

– $G2$: $S \rightarrow I \mid IT$
 $T \rightarrow I \mid d \mid IT \mid dT$

正规文法和正规式

对任意一个正规文法，存在一个定义同一个语言的正规式；反之亦然。

■ Σ 上的正规式 \Rightarrow 正规文法

■ 初始 $V_T = \Sigma, S \in V_N$ ，生成正规产生式(或定义式) : $S \rightarrow r$
(r 为正规式)

■ (R1) 对形如 $A \rightarrow r_1 r_2$ 的正规产生式:

$A \rightarrow r_1 B$

$B \rightarrow r_2$

$B \in V_N$

(R2) 对形如 $A \rightarrow r^* r_1$ 的正规产生式:

$A \rightarrow r B$

$A \rightarrow r_1$

$B \rightarrow r B$

$B \rightarrow r_1 \quad B \in V_N$

(R3) 对形如 $A \rightarrow r_1 \mid r_2$ 的正规产生式:

$A \rightarrow r_1$

$A \rightarrow r_2$

■ 不断应用R做变换，直到每个产生式右端只含一个 V_N

正规文法和正规式

- 例 $r=a(a \mid d)^*$ 转换为正规文法

$V_T=\{a,d\}$ $S \rightarrow a(a \mid d)^*$

– R1 $S \rightarrow aA$

$A \rightarrow (a \mid d)^*$

– R2 $A \rightarrow (a \mid d)B$

$A \rightarrow \varepsilon$

$B \rightarrow (a \mid d)B$

$B \rightarrow \varepsilon$

正规文法和正规式

– R3 $G[s]:$

P {

- $S \rightarrow aA$
- $A \rightarrow \varepsilon$
- $A \rightarrow aB$
- $A \rightarrow dB$
- $B \rightarrow aB$
- $B \rightarrow dB$
- $B \rightarrow \varepsilon$

$V_T = \{a, d\}$

$V_N = \{S, A, B\}$

正规文法和正规式

- 例：将标示符集合的正规式：
 $\text{letter}(\text{letter}|\text{digit})^*$ 转换为正规文法

$S \rightarrow \text{letter}A$

$A \rightarrow (\text{letter}|\text{digit})^*$

$\left\{ \begin{array}{l} A \rightarrow (\text{letter}|\text{digit})B | \varepsilon \end{array} \right.$

$\left\{ \begin{array}{l} B \rightarrow (\text{letter}|\text{digit})B | \varepsilon \end{array} \right.$

- 结果：
$$\left\{ \begin{array}{l} S \rightarrow \text{letter}A \\ A \rightarrow \text{letter } B | \text{digit}B | \varepsilon \\ B \rightarrow \text{letter}B | \text{digit}B | \varepsilon \end{array} \right.$$

- 例：整数的正规式 $(\text{digit})^+$

- 运算符的正规式 $\text{relop} \rightarrow < | <= | = | <> | > | >=$

正规文法和正规式

正规文法 \Rightarrow 正规式

正规文法

正规式

- $A \rightarrow xB, B \rightarrow y$ 转换成: $A = xy$
- $A \rightarrow xA \mid y$ 转换成: $A = x^*y$
- $A \rightarrow x \mid y$ 转换成: $A = x \mid y$

正规文法和正规式

– $G[s]: S \rightarrow aA$

$S \rightarrow a$

$A \rightarrow aA$

$A \rightarrow dA$

$A \rightarrow a$

$A \rightarrow d$

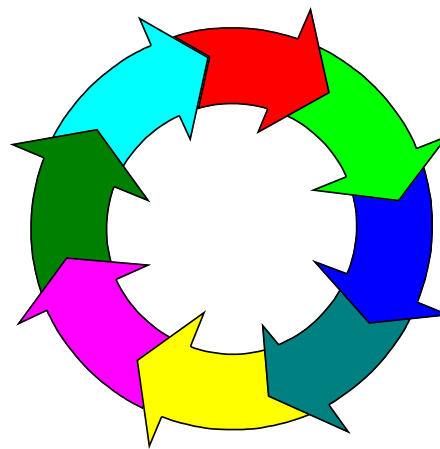
- ① $S \rightarrow aA \mid a$
② $A \rightarrow aA \mid a \mid dA \mid d$
③ $A \rightarrow (a \mid d)A \mid (a \mid d)$
④ $A \rightarrow (a \mid d)^*(a \mid d)$

正规文法和正规式

$$\begin{aligned} - S &= a(a \mid d)^*(a \mid d) \mid a \\ &= a((a \mid d)^*(a \mid d) \mid \varepsilon) \\ &= a((a \mid d)^+ \mid \varepsilon) \\ - R &= a(a \mid d)^* \end{aligned}$$

4.3有穷(限)自动机

- 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



4.3.1 DFA

■ DFA定义:

- 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、 Σ 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 Σ 为输入符号字母表；

DFA定义

- 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。

DFA 例:

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

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

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

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

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

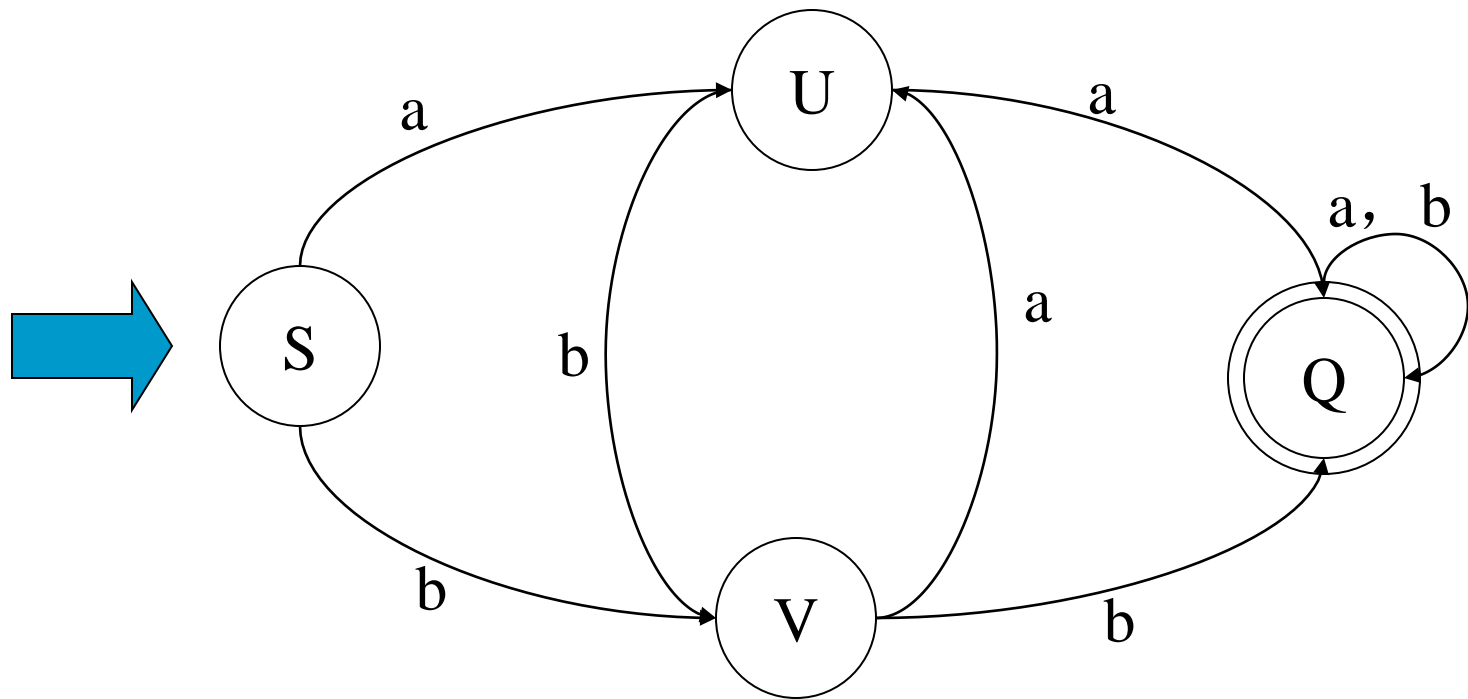
$$- f(U, a) = Q$$

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

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

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

DFA 的状态转换图表示



DFA 的矩阵表示

状态 \ 字符	a	b
S	U	V
U	Q	V
V	U	Q
Q	Q	Q

DFA

■ DFA M的作用：

对于 Σ^* 中的任何字符串 t ，若存在一条从初态结到某一终态结的道路，且这条路上所有弧的标记符连接成的字符串等于 t ，则称 t 可为**DFA M**所接受（识别）。若**M**的初态结同时又是终态结，则空字可为**M**所识别（接受）。

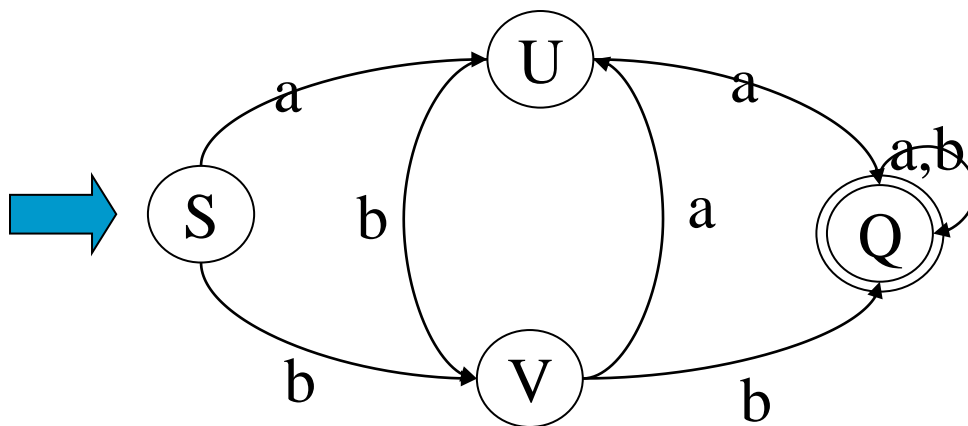
DFA

- Σ^* 上的符号串 t 被 M 接受的形式叙述：
 - 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）。
- Σ^* 上的符号串 t 在 M 上运行的定义：
 - 一个输入符号串 t , （我们将它表示成 $t_1 t_x$ 的形式, 其中 $t_1 \in \Sigma$, $t_x \in \Sigma^*$ ）在DFA M 上运行的定义为：
 - $f(Q, t_1 t_x) = f(f(Q, t_1), t_x)$ 其中 $Q \in K$

DFA

■ 例：证明 $t=baab$ 被前例中的DFA所接受。

- $f(S, baab) = f(f(S, b), aab) = f(V, aab)$
 $= f(f(V, a), ab) = f(U, ab) = f(f(U, a), b)$
 $= f(Q, b) = Q$
- Q属于终态。
- 得证。

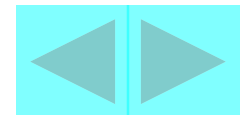


DFA

- DFA M 所能接受的符号串的全体记为 $L(M)$ (语言)
- 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 $M = (K, \Sigma, f, S, Z)$ 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')



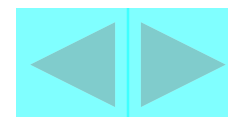
4.3.2不确定的有穷自动机NFA

■ 定义

- $N = \{K, \Sigma, f, S, Z\}$, 其中 K 为状态的有穷非空集, Σ 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, $S \subset K$ 是初始状态集, $Z \subset K$ 为终止状态集。

■ 例子

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



4.3.2 NFA

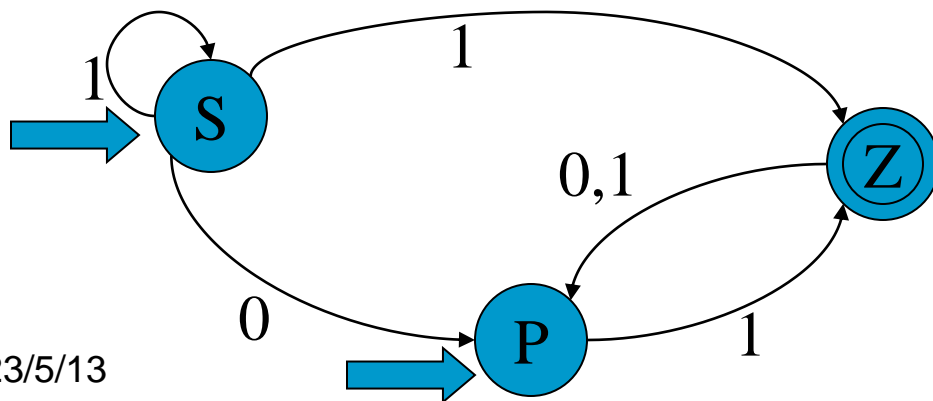
– $f(S, 1) = \{S, Z\}$

– $f(P, 1) = \{Z\}$

– $f(Z, 0) = \{P\}$

– $f(Z, 1) = \{P\}$

■ 状态图表示



4.3.2 NFA

■ 例： 一个NFA $M =$

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

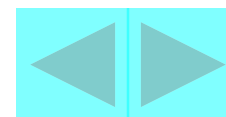
$f(0,a) = \{0,3\}$ $f(0,b) = \{0,1\}$

$f(1,b) = \{2\}$ $f(2,a) = \{2\}$

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

$f(4,a) = \{4\}$ $f(4,b) = \{4\}$

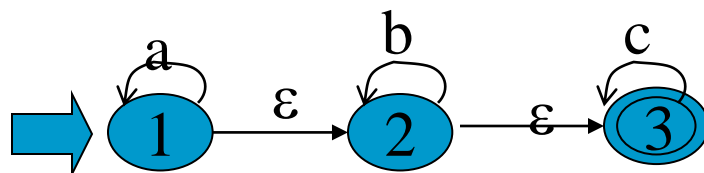
■ 画出其状态图



4.3.2 NFA

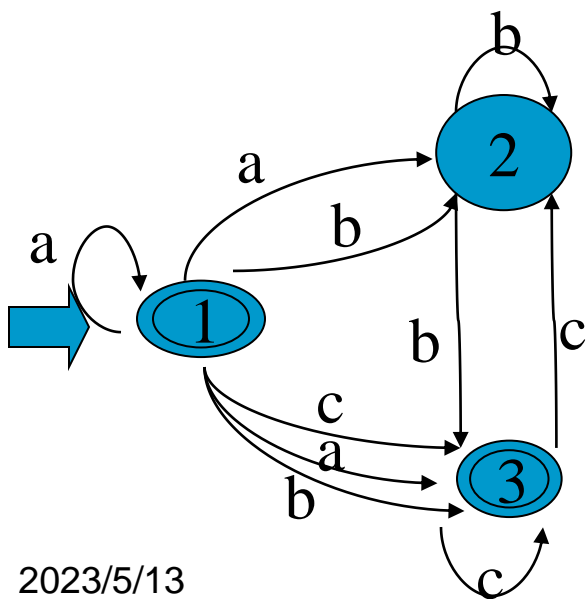
- Σ^* 上的符号串 t 在NFA N 上运行...
- Σ^* 上的符号串 t 被NFA N 接受...
- 具有 ε 转移的不确定的有穷自动机NFA... f 为 $K \times (\Sigma \cup \{\varepsilon\})$ 到 K 的子集 (2^K) 的一种映射

—



4.3.2 NFA

- 对任何一个具有 ϵ 转移的不确定的有穷自动机NFA N ，一定存在一个不具有 ϵ 转移的不确定的有穷自动机NFA M ，使得 $L(M)=L(N)$ 。



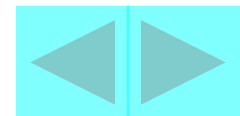
4.3.3 NFA的确定化

- 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 。
与某一NFA等价的DFA不唯一



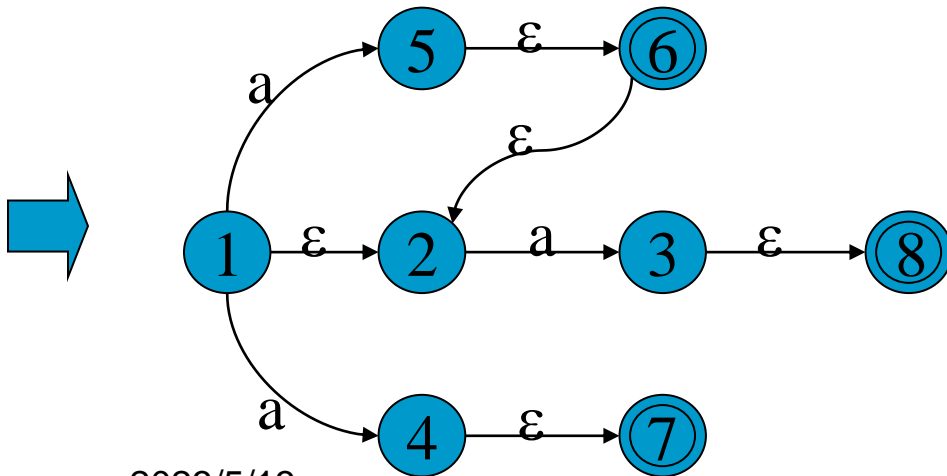
4.3.3 NFA的确定化

- 定义对状态集合 I 的几个有关运算：
- 1、状态集合 I 的 ε -闭包，表示为 $\varepsilon\text{-closure}(I)$ ，定义为一状态集，是状态集 I 中的任何状态 S 经任意条 ε 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 $\varepsilon\text{-closure}(I)$ 。
- 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 的某一状态经过一条 a 弧而到达的状态的全体。定义 $Ia = \varepsilon\text{-closure}(J)$



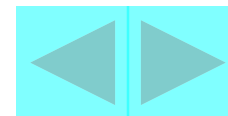
4.3.3 NFA的确定化

- $I = \{1\}$, $\varepsilon\text{-closure}(I) = \{1, 2\}$;
- $I = \{5\}$, $\varepsilon\text{-closure}(I) = \{5, 6, 2\}$;
- $\text{move}(\{1, 2\}, a) = \{5, 3, 4\}$
- $\varepsilon\text{-closure}(\{5, 3, 4\}) = \{2, 3, 4, 5, 6, 7, 8\}$;



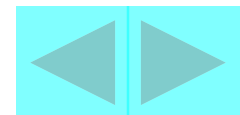
4.3.3 NFA的确定化

- 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 $[S_1 S_2]$;



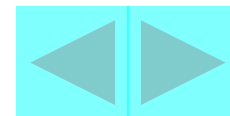
4.3.3 NFA的确定化

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



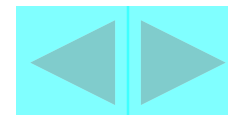
4.3.3 NFA的确定化

- 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 $C = (T_1, T_2, \dots, T_l)$ ，其中 T_1, T_2, \dots, T_l 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。



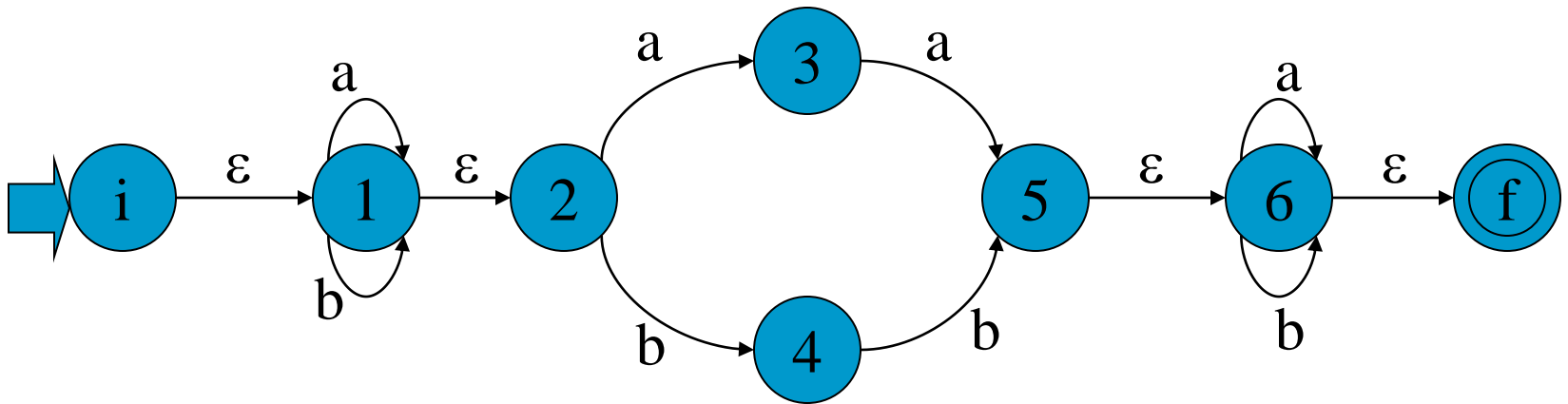
4.3.3 NFA的确定化

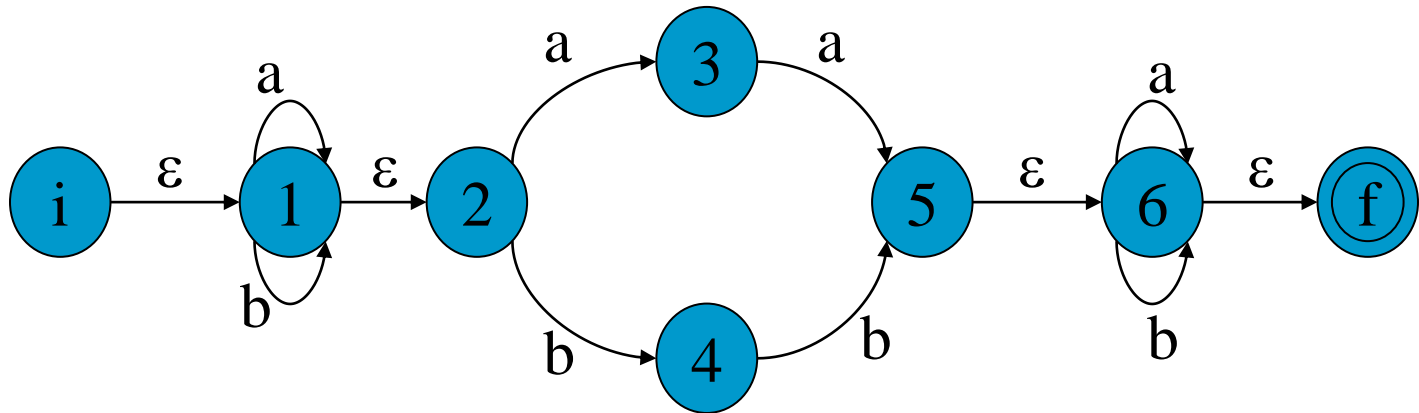
```
– 2  while (C中存在尚未被标记的子集T)
      do  {
            标记T;
            for 每个输入字母a  do
            {
                  U:=  $\varepsilon$ -closure(move(T,a));
                  if U不在C中  then
                        将U作为未标记的子集加
                        在C中
            }
      }
```



4.3.3 NFA的确定化

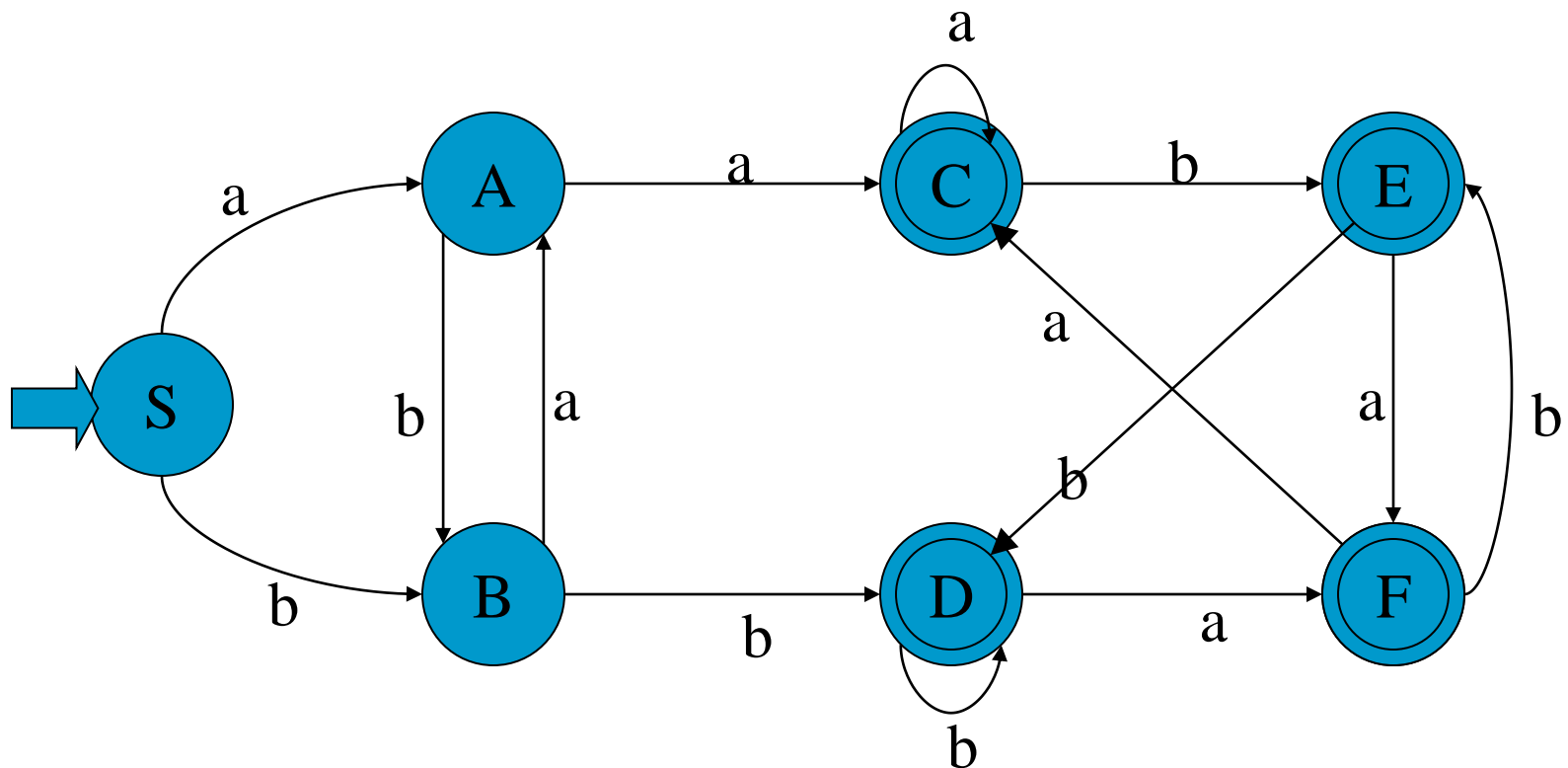
■ 例子





		Ia	Ib
{i,1,2}	S	{1,2,3}	A
{1,2,3}	A	{1,2,3,5,6,f}	C
{1,2,4}	B	{1,2,3}	A
{1,2,3,5,6,f}	C	{1,2,3,5,6,f}	C
{1,2,4,5,6,f}	D	{1,2,3,6,f}	F
{1,2,4,6,f}	E	{1,2,3,6,f}	F
{1,2,3,6,f}	F	{1,2,3,5,6,f}	C

4.3.3 NFA的确定化



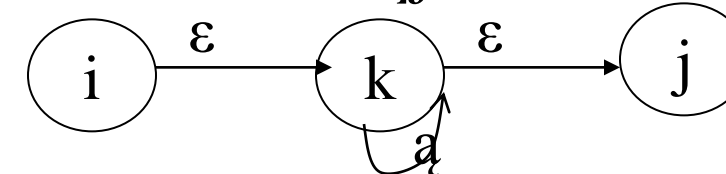
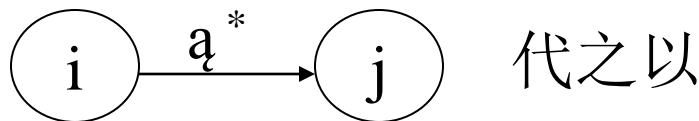
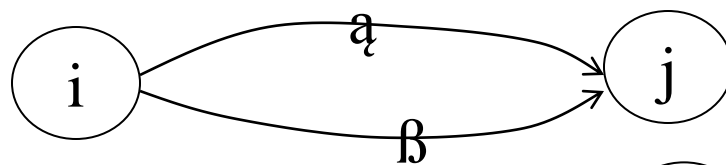
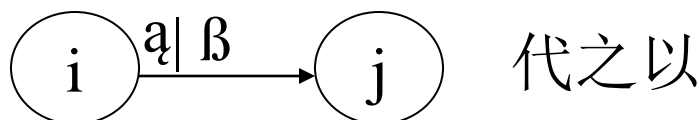
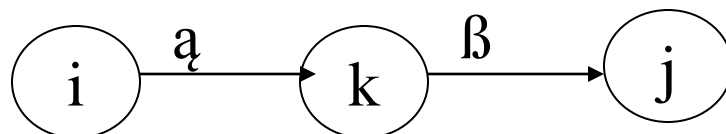
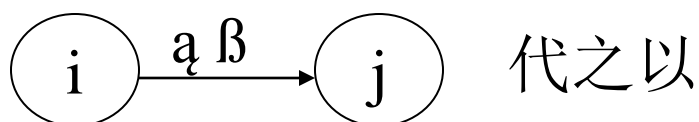
4.4有穷自动机和正规表达式

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 $L(M)=L(R)$ 。

正规式 \Rightarrow 有穷自动机

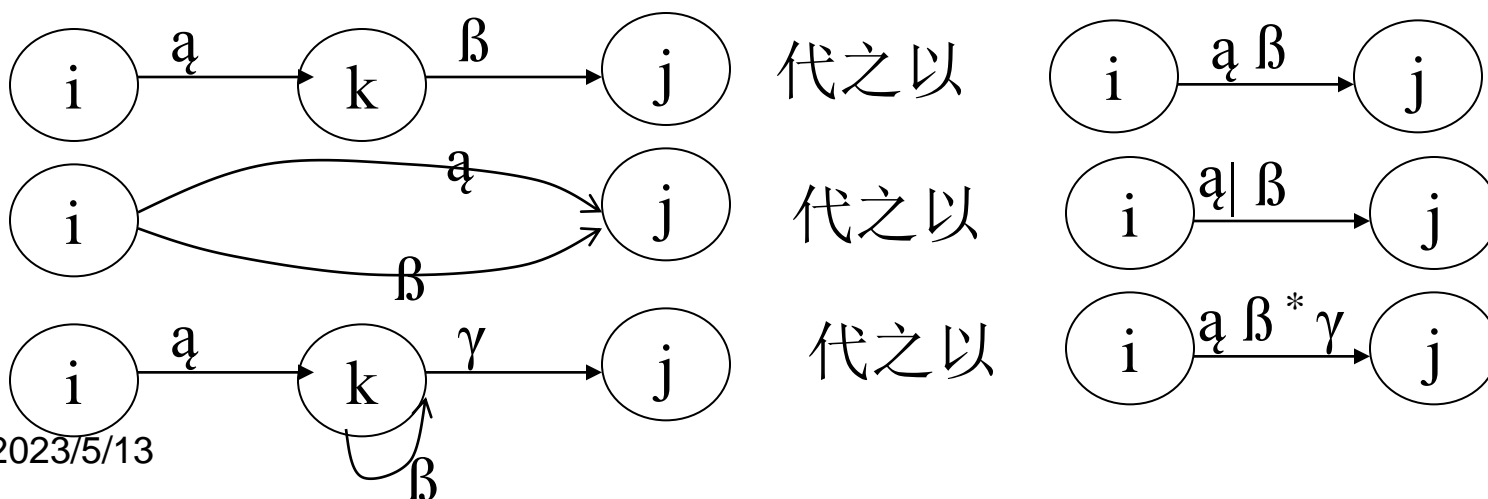
- 设给定正规式 W ，则构造相应自动机的方法如下：
- 若 $W=\emptyset$ ，则对应的NFA为 $\longrightarrow (S) \quad (T)$
- 若 $W \neq \emptyset$ ，则对应的NFA为 $\longrightarrow (S) \xrightarrow{W} (T)$
- 然后利用以下规则加入结点和箭弧，直到得到自动机为止。



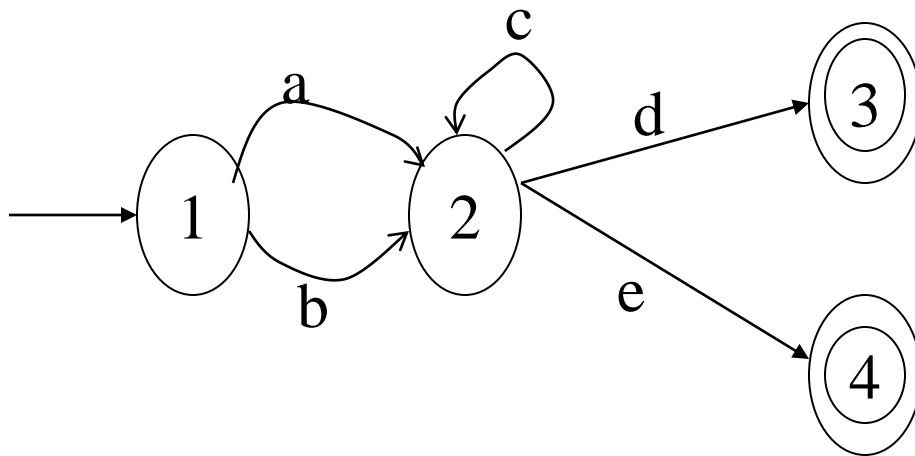
- 例：正规式 $(a|b)^*(aa|bb)(a|b)^*$ 转化为自动机

有穷自动机=>正规式

- 首先检查是否只有一个终态结点，若有多个，则引入新结点T，从所有终态结点引 ϵ 边到T结点，并令T为唯一的终态结点。
- 然后按以下规则消除结点与箭弧：

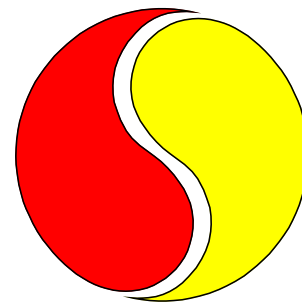


■ 例:



4.5有穷自动机和正规文法

- 有穷自动机和正规文法的等价性:
 - 1.对于一个NFA M，都存在一个正规文法G，使得 $L(G)=L(M)$.
 - 2.对于一个正规文法G，都存在一个NFA M，使得 $L(M)=L(G)$.



正规文法=>自动机

- 字母表与G的终结符集相同;
- G中的非终结符对应状态, 开始符对应开始状态
- 增加一个新的终结状态Z。
- G中的 $A \rightarrow tB$ 构造转换函数 $f(A,t)=B$
- G中的 $A \rightarrow t$ 构造转换函数 $f(A,t)=Z$

■ 例：求与G[S]等价的NFA:

■ G[S]:

■ $S \rightarrow aA$

■ $S \rightarrow bB$

■ $S \rightarrow \varepsilon$

■ $A \rightarrow aB$

■ $A \rightarrow bA$

■ $B \rightarrow aS$

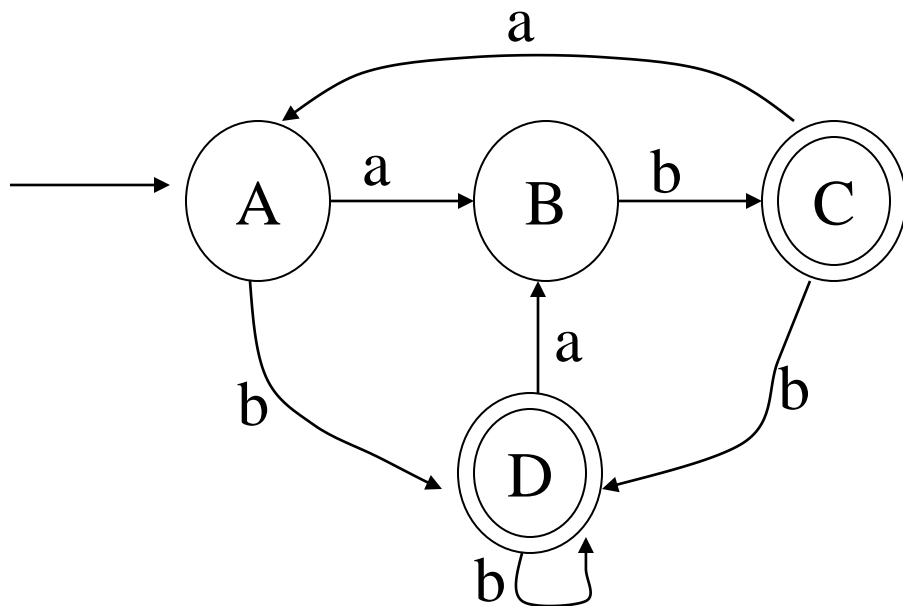
■ $B \rightarrow bA$

■ $B \rightarrow \varepsilon$

自动机=>正规文法

- 对转换函数 $f(A,t)=B$ ，对应产生式：
 $A \rightarrow tB$
- 对终态 Z ，增加一产生式： $Z \rightarrow \varepsilon$
- NFA的初态对应文法的开始符号；
- NFA的字母表对应文法的终结符号集。

■ 例：给出与下图NFA等价的正规文法G。



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

$A \rightarrow aB$

$C \rightarrow \epsilon$

$A \rightarrow bD$

$D \rightarrow aB$

$B \rightarrow bC$

$D \rightarrow bD$

$C \rightarrow aA$

$D \rightarrow \epsilon$

总结

- 词法分析程序
- 正规表达式与正规集（正规语言）
- 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