

# 编译原理

## Compiler Construction Principles



朱 青

信息学院计算机系，  
中国人民大学，

zqruc2012@aliyun.com



# 第2章:词法分析 (Lexical Analysis)

---

⌘ 2.1 词法分析程序的功能

⌘ 2.2 词法分析器的设计

⌘ 2.3 正规表达式 (Regular Expression)

⌘ 2.4 有限自动机

⌘ 2.5 词法分析器的自动生成

## 2.4 有限自动机

---

⌘ 2.4.1 确定有限自动机 (DFA)

⌘ 2.4.2 非确定有限自动机 (NFA) 确定化

⌘ 2.4.3 具有 $\varepsilon$ -转移的NFA M确定化

⌘ 2.4.4 DFA的化简

⌘ 2.4.5 正规式与有限自动机的等价性

⌘ 2.4.6 正规文法与有限自动机

## 2.4.5 正规式与有限自动机的等价性

---

定义1：（子集法）若 $I$ 为状态集 $S$ 的子集。

$$Ia = \varepsilon\text{-closure}(J) \quad a \in \Sigma$$

则称： $J$ 是从 $I$ 出发，经过一条 $a$ 弧所能到达的状态子集。

---

定理1:

$\Sigma$ 上的NFA  $M$  所识别的字的全体构成 $\Sigma$ 上的正规式 $V$ .

定理2:

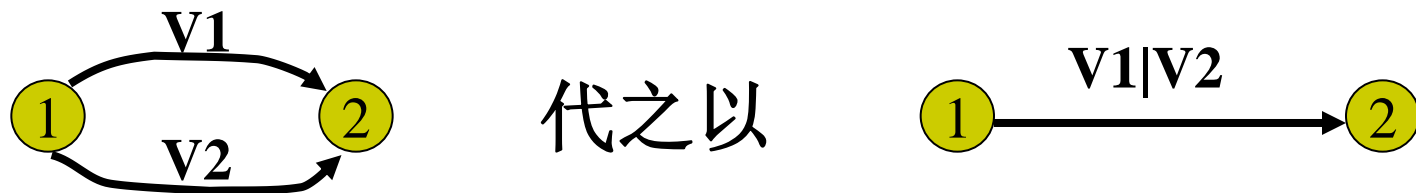
对 $\Sigma$ 上的任何正规集,总存在一个 DFA  $M$ ,使 $L(V)=L(M)$ .

证明1:

把转换图的意义拓宽,令每条弧上可以标记正规式.

在给定的NFA  $M$ 上加二个新结,  
一个为初态 $X$ ,从 $X$ 用 $\epsilon$ 弧连接 $M$ 的所有  
初态,另一为 $Y$ ,从 $M$ 的所有终态用 $\epsilon$ 弧连到 $Y$ ,新  
的NFA  $M'$ 与 $M$ 等价.

对 $M'$ 用以下等价规则,消去除 $X,Y$ 以外的其它结点.

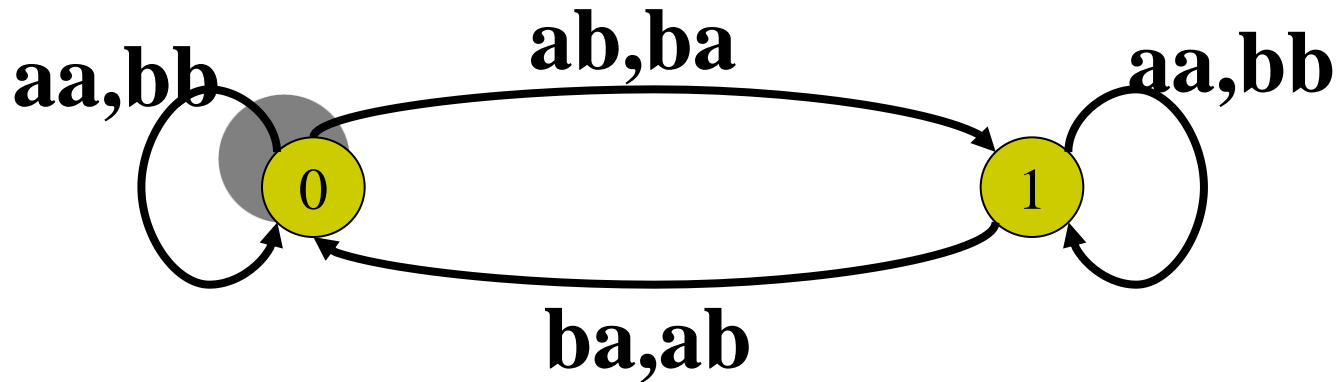


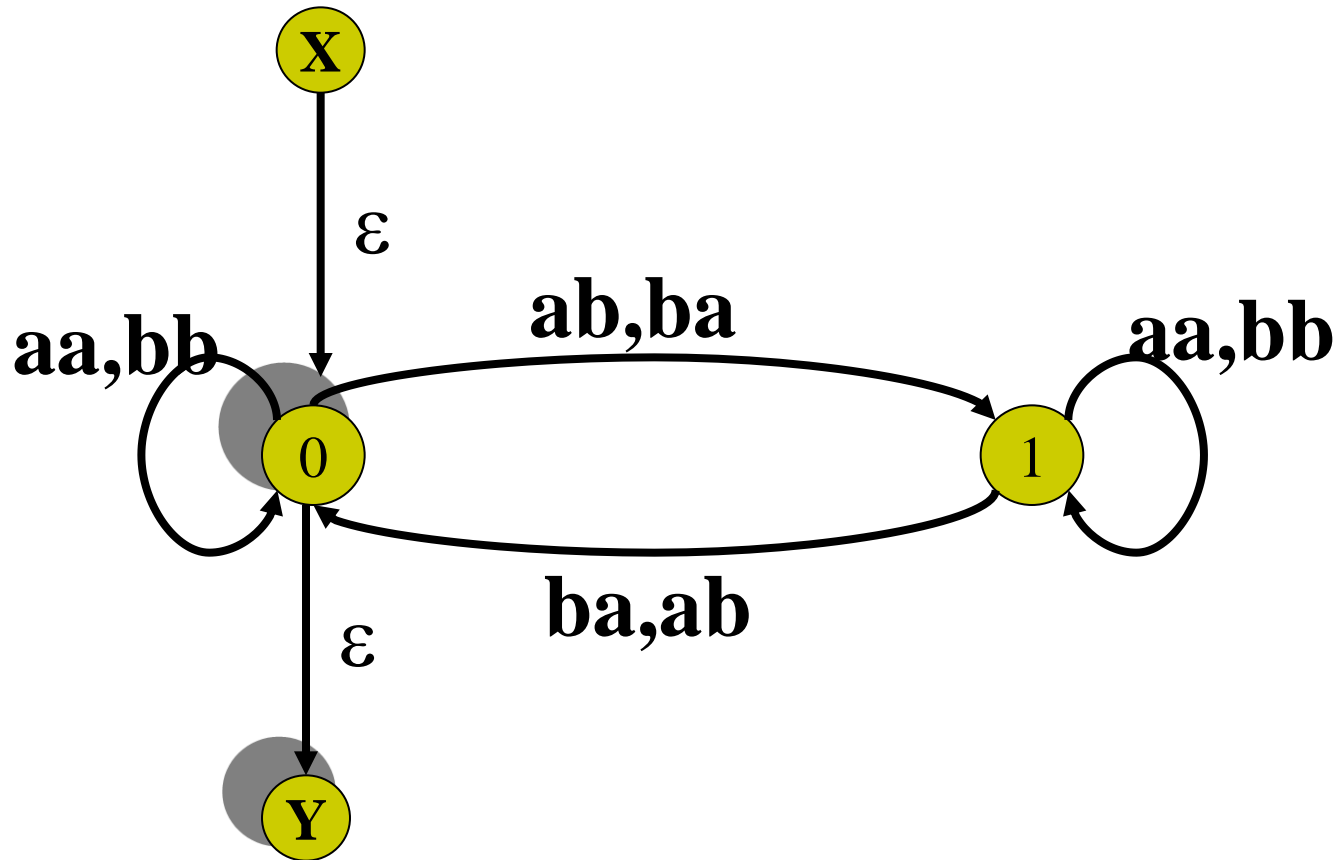
于是,最后得到只有初态X和终态Y的NFA,其标记为一正规式.即:

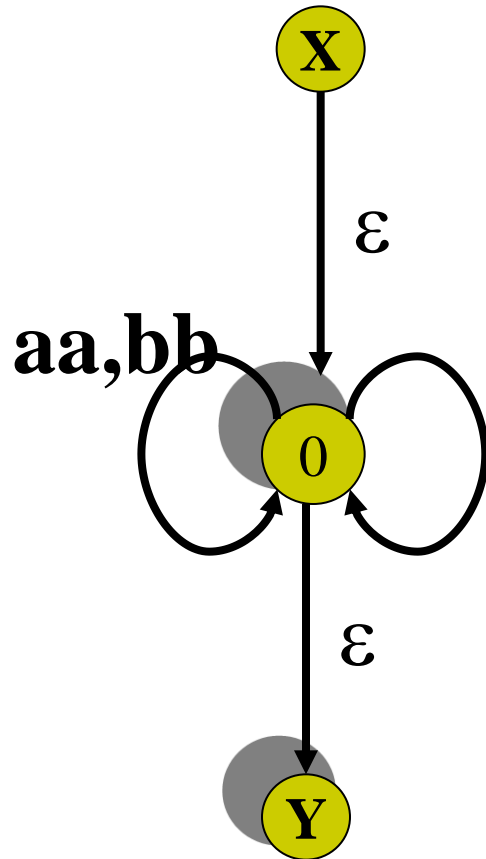




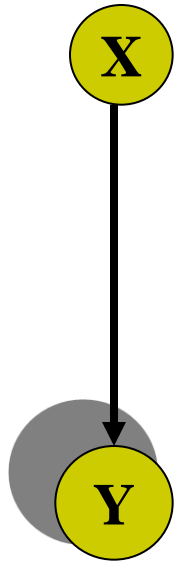
例如:把下面的NFA  $M \Rightarrow$  正规式.  
NFA  $M$  是识别具有偶数a和偶数个b的非有限自动机:







$(ab|ba)(aa|bb)^*(ba|ab)$



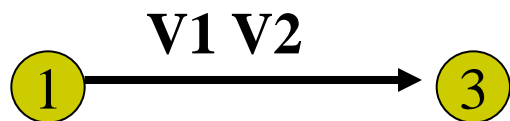
**$((aa|bb)|(ab|ba)(aa|bb)^*(ba|ab))^*$**

证明2：分两步：

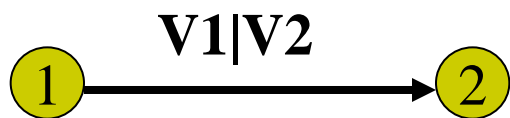
1) 对给定的正规式构成一个NFA M。  
先写出：



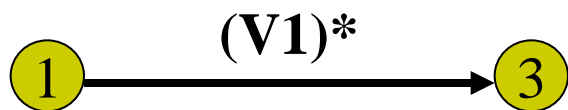
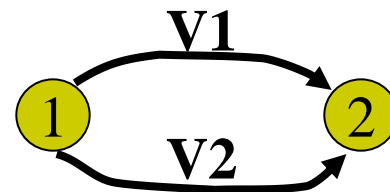
用以下规则对V进行分解并加进新结。



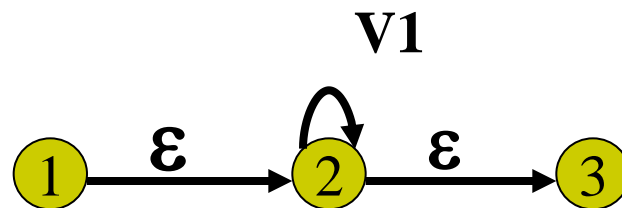
代之以



代之以



代之以



在分解过程中，要求：

- (1)  $X$ ,  $Y$ 始终为唯一的初态和终态。
- (2) 所加新结其名字彼此不同。
- (3) 弧上的标记必须是字符或空字。

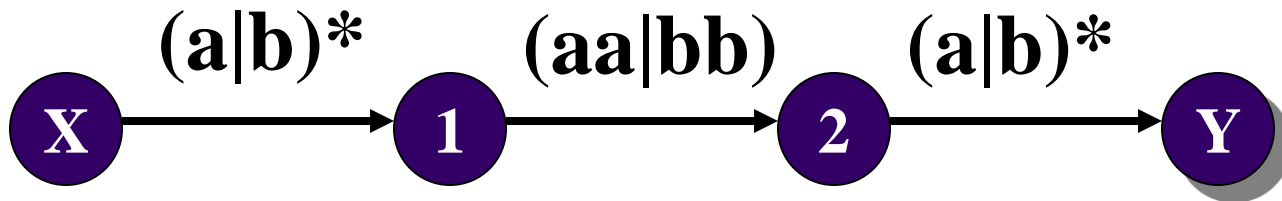
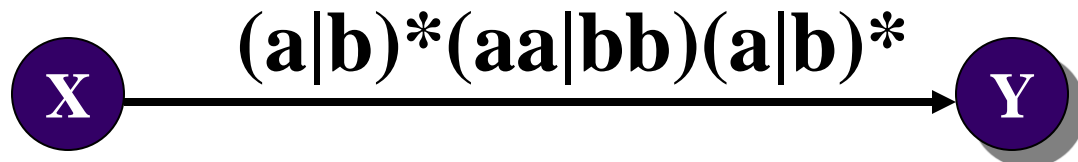
2) 把NFA M 确定化：采用子集法进行确定化。  
为简单记，设： $\Sigma=\{a, b\}$ .  
因为一个状态矩阵可唯一的刻画一个DFA，故  
用子集法造该矩阵。

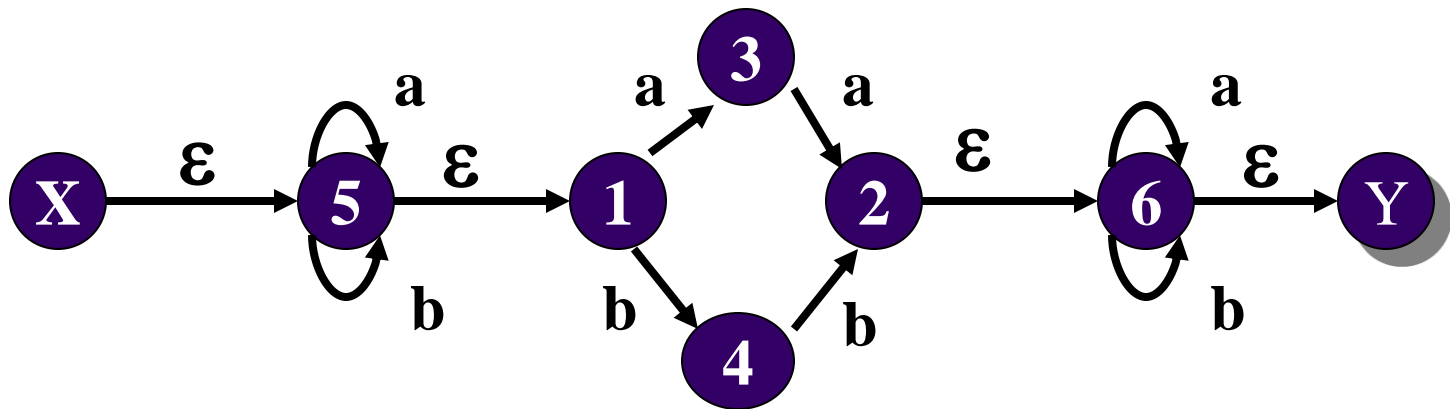
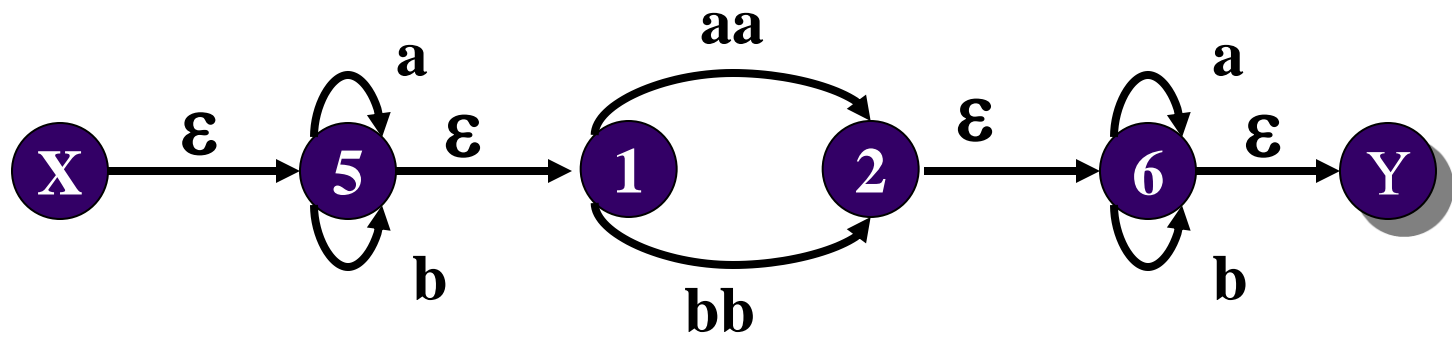
状态子集	Ia	Ib
------	----	----

(在前面章节中已讲)

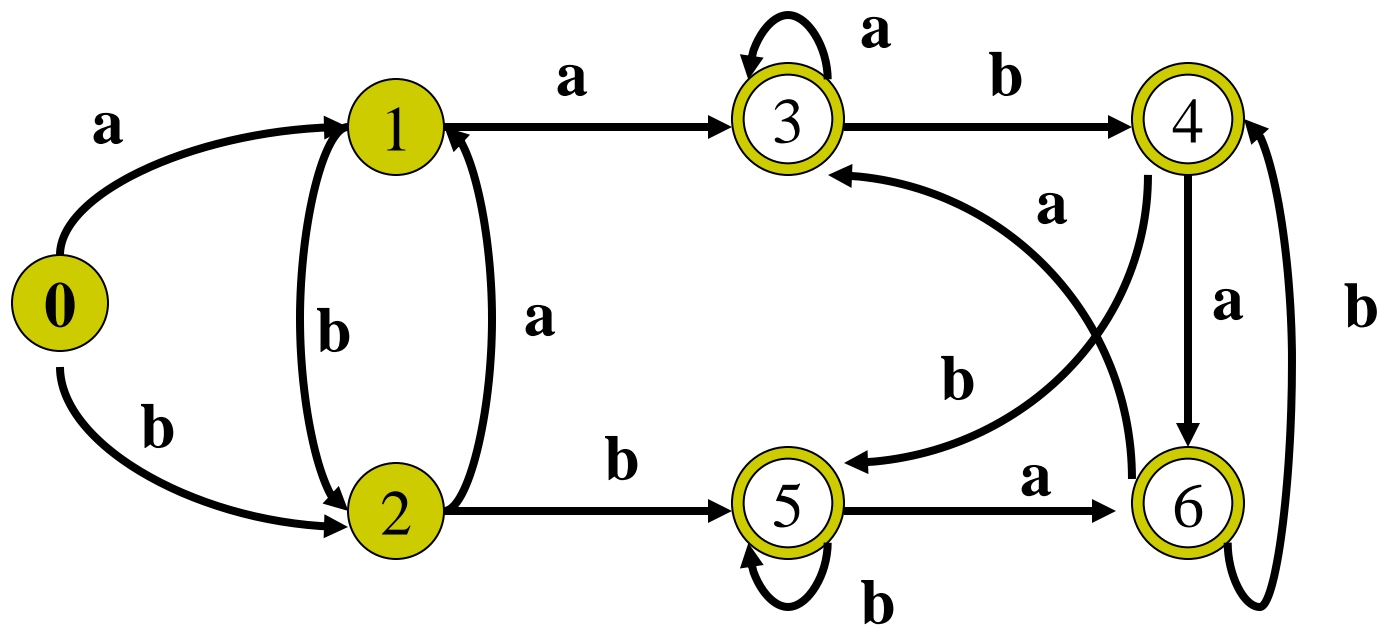


例： 设  $V=(a|b)^*(aa|bb)(a|b)^*$





I	Ia	Ib
0{X,5,1}	{5,3,1}	{5,4,1}
1{5,3,1}	{5,3,1,2,6,Y}	{5,4,1}
2{5,4,1}	{5,3,1}	{5,3,1,2,6,Y}
3{5,3,1,2,6,Y}	{5,3,1,2,6,Y}	{5,4,1,6,Y}
4{5,4,1,6,Y}	{5,3,1,6,Y}	{5,3,1,2,6,Y}
5{5,4,1,2,6,Y}	{5,3,1,6,Y}	{5,4,1,2,6,Y}
6{5,3,1,6,Y}	{5,3,1,2,6,Y}	{5,4,1,6,Y}

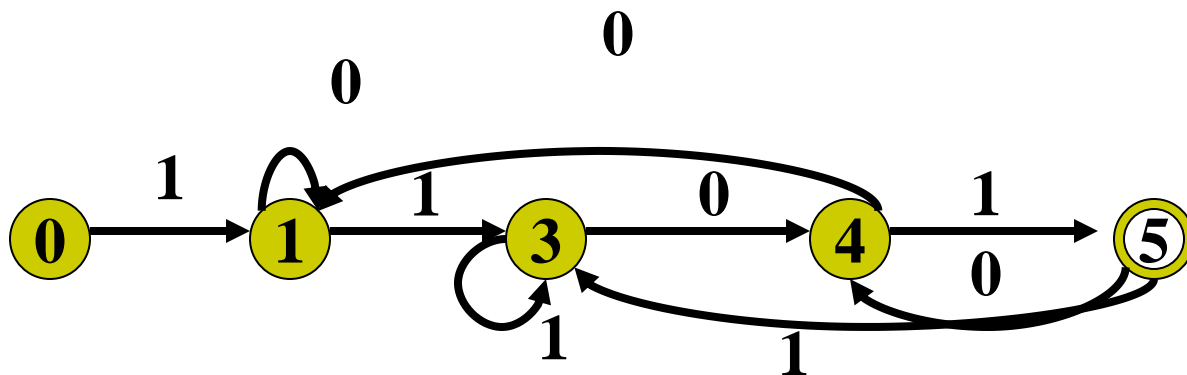


推论：

一个子集是正规的，当且仅当  
他可由一个**DFA**（或**NFA**）所  
识别。

例题：构造下列正规式相应的**DFA**  
 $1(0|1)^*101$

答案：



## 2.4.6 正规文法与有限自动机

正规文法是描述单词符号的  
另一种方法.

例如:

标识符  $\longrightarrow$  字母 |  $\langle$ 标识符 $\rangle$  字母  
                                  |  $\langle$ 标识符 $\rangle$  数字.

<整常数>  $\longrightarrow$  数字 | <整常数> 数字

<分界符>  $\longrightarrow$  + | - | \* | / | , | ; | ( | ) | ...

<分界符>  $\longrightarrow$  <冒号> = | <星号> \* | <斜竖> / | ...

<冒号>  $\longrightarrow$  :

<星号>  $\longrightarrow$  \*

<斜竖>  $\longrightarrow$  /

定义式有两种类型:  $P \longrightarrow t$

$P \longrightarrow Qt$



# (1)正规文法

定义:

1) 如果文法 $G=(V_T, V_N, S, P)$

其中: (1)  $V_T$ 是非空有限集, 每个元素是一个终结符.

(2)  $V_N$ 是非空有限集, 每个元素是一个非终结符.

- (3)  $S$ 是一个非终结符,是开始符号.  
( $S$ 在产生式的左部必须至少出现一次)
- (4)  $P$ 是产生式的集合:它的每一个产生式 $P$ 的形式为:

$$A \longrightarrow aB \quad \text{或} \quad A \longrightarrow a$$

其中,  $A, B \in V_N$ ,  $a \in V_T \cup \{\epsilon\}$ ,

则称  $G$  是右线性文法。

2) 若文法G中的每一个产生式的形式为

$$A \longrightarrow Ba \text{ 或 } A \longrightarrow a$$

则称 G是左线性文法。

3) 右线性文法和左线性文法都称为3型文法. 3型文法也称为正规文法. 它所产生的语言称为3型语言或正规语言.

## (2)正规文法构造相应的

### 状态转换图:

1) 对于右线性文法:

设: $G=\{V_N, V_T, P, S\}$ 是一个右线性文法, 并设 $|V_N|=k$ ,

构造的状态转换图中共有 $k + 1$ 个结点.

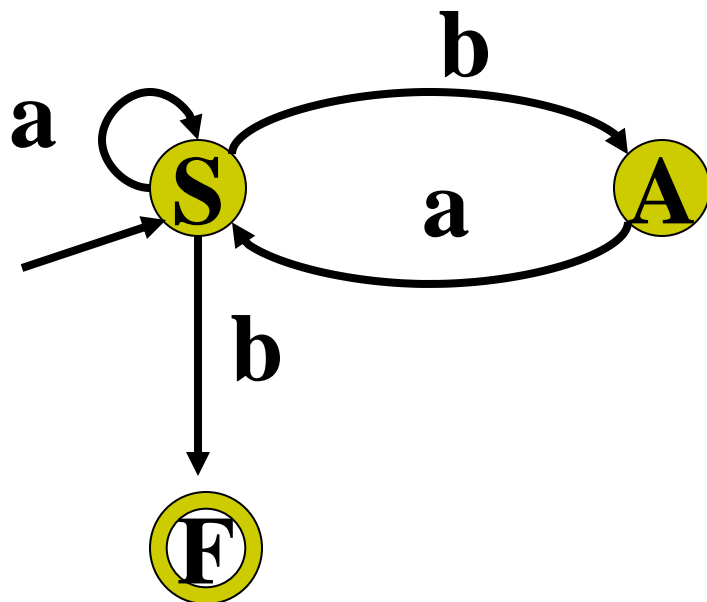
结点的标记: ----- 用 $V_N$   
的各个非终结符号分别标记其  
中的 $k$  个结点, 且令 $G$ 的开始符  
号 $S$ 作为初态, 余下的一个作为  
终态结点 $F$ , 且 $F$ 不属于 $V_N$ .

## 箭弧的规则:

对于**G**中每一形如 $A \xrightarrow{a} B$  的产生式, 从结点**A**引一条箭弧到结点**B**,并用符号**a**标记这条箭弧.

对于**G**中每一形如 $A \xrightarrow{a}$ 的产生式,从结点**A**引一条箭弧到终态结点**F**,并用符号**a**标记这条箭弧.

例如： 设给定右线性文法G：

$$S \longrightarrow aS \mid bA \mid b$$
$$A \longrightarrow aS$$


2) 对于左线性文法:

设: $G=\{V_N, V_T, P, S\}$ 是一个左线性文法, 并设 $|V_N|=k$ ,

构造的状态转换图 $M$ 中共有 $k+1$ 个结点.



结点的标记: ----- 用 $V_N$   
的各个非终结符号分别标记其  
中的 $k$ 个结点, 且 引入 开始符  
号 $R$  ( $R$ 不属于 $V_N$ ) 作为初态,  
用 $S$ 作为终态结点.

## 箭弧的规则:

对于**G**中每一形如**A**  $\rightarrow$  **Ba** 的产生式, 从结点**B**引一条箭弧到结点**A**,并用符号**a**标记这条箭弧.

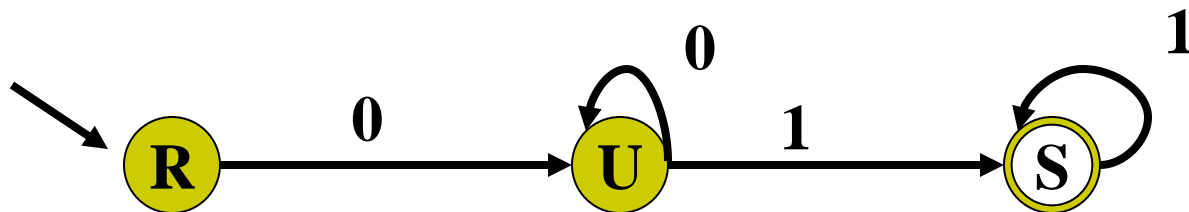
对于**G**中每一形如**A**  $\rightarrow$  **a**的产生式,从初态**R**引一条箭弧到结点**A**,并用符号**a**标记这条箭弧.

例如：对于左线性文法  $G = (\{S, U\}, \{0, 1\}, P, S)$  其中

$$P = \{ S \longrightarrow S1, S \longrightarrow U1, \\ U \longrightarrow U0, U \longrightarrow 0 \}$$

构造状态转换图。

解：



定理：对任何一个右线性正规文法  $G$ ，都存在一个左线性正规文法  $G'$ ，  
使  $L(G') = L(G)$ ，  
反之亦然。

# (3) 正规文法与有限自动机的等价性

定理1：对每一个 右线性正规文法G或左线性正规文法G，都存在有限自动机M，使 $L(M) = L(G)$ 。

定理2: 对每一个DFA  $M$ , 都存在  
一个右线性正规文法  $G$  和  
一个左线性正规文法  $G'$ ,  
使  
 $L(M) = L(G') = L(G)$ 。

由DFA M定义右线性文法:

设: DFA  $M = (\Sigma, V_N, S, F, f)$ ,

如初态符号  $S$  不属于  $F$ , 令

$$G = (\Sigma, V_N, S, P),$$

其中:  $\Sigma$  为终结符.  $P$  是按下面规则定义: 对

任何  $a$  属于  $\Sigma$  且  $A, B$  属于  $V_N$ , 若

有  $f(A, a) = B$  则

a) 当  $B$  不属于  $F$ , 令  $A \longrightarrow aB$

b) 当  $B$  属于  $F$ , 令  $A \longrightarrow a \mid aB$

- 如初态符号 **S** 属于 **F**, 因为  $f(S, \varepsilon) = S$ , 所以,  $\varepsilon$  属于  $L(M)$ , 但  $\varepsilon$  不属于上面构造的 **G** 所产生的语言  $L(G)$ . 实际上

$$L(G) = L(M) - \{\varepsilon\}$$

因而对上面由 **M** 出发所构造的右线性正规文法 **G** 中添加一个非终结符号  $S_0$  不属于  $V_N$  和产生式

$$S_0 \longrightarrow S \mid \varepsilon$$



并用 $S_0$ 代替 $S$ 作开始符号.这样经过修正后的 $G$ 仍是右线性正规文法且 $L(G)=L(M)$ .

类似的从 $M$ 出发可构造左线性文法.

例题： 设DFA  $M=(\{0,1\},\{A,B,C,D\}$   
A,{B},f )

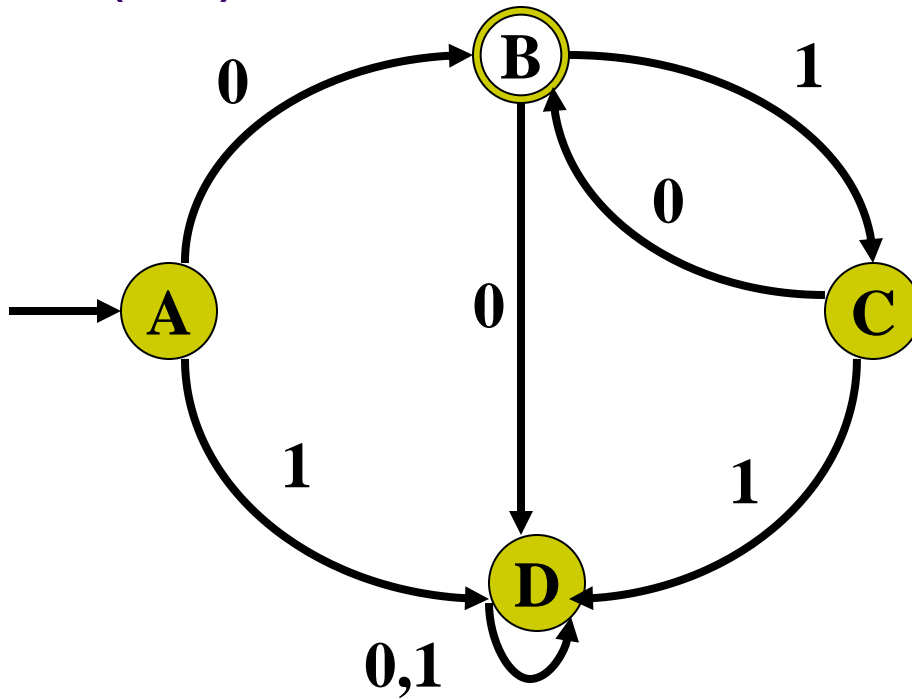
$$f(A,0)=B \quad f(A,1)=D$$

$$f(B,0)=D \quad f(B,1)=C$$

$$f(C,0)=B \quad f(C,1)=D$$

$$f(D,0)=D \quad f(D,1)=D$$

$$L(M)=0(10)^*.$$



构造右线性正规文法：

$$G = (\{0, 1\}, \{A, B, C, D\}, A, P)$$

其中P为产生式的集合：

$$A \longrightarrow 0|0B|1D$$

$$B \longrightarrow 0D|1C$$

$$C \longrightarrow 0|0B|1D$$

$$D \longrightarrow 0D|1D$$

$$L(G) = L(G) = 0(10)^*$$

由右线性正规文法**G**出发构造的NFA **M'**  
为

$$M' = (\{0, 1\}, \{A, B, C, D, f\}, A, \{f\}, F)$$

其中

$$F(A, 0) = \{f, B\}$$

$$F(A, 1) = \{D\}$$

$$F(B, 0) = \{D\}$$

$$F(B, 1) = \{C\}$$

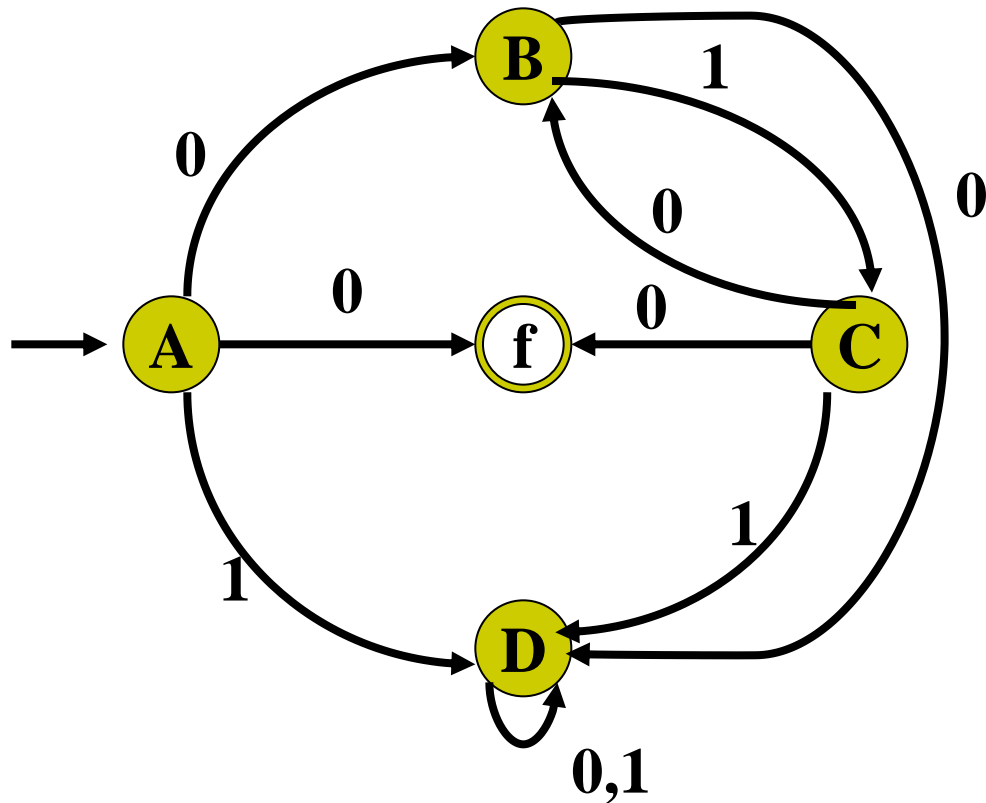
$$F(C, 0) = \{f, B\}$$

$$F(C, 1) = \{D\}$$

$$F(D, 0) = \{D\}$$

$$F(D, 1) = \{D\}$$

状态转换图如下所示:



$$L(M') = L(M)$$

据NFA  $M'$ 的状态转换图，构造左线性正规文法

$G' = (\{0, 1\}, \{B, C, D, S\}, S, P)$

其中 $S=f$  为下面产生式的集合:

$S \longrightarrow 0 \mid C0 \quad C \longrightarrow B1$

$B \longrightarrow 0 \mid C0$

$D \longrightarrow 1 \mid C1 \mid D0 \mid D1 \mid B0$

这里从 $M'$ 构造左线性文法 $G'$ 的产生式 $P$ 的方法:

若对任何 $A, A1$ 属于 $\{B, C, D, S\}$   
且 $a$ 属于 $\{0, 1\}$ 有 $F(A1, a) = A2$ ,

则令 $A2 \longrightarrow A1a$ ,

若有 $F(A1, a) = A$  ( $A$ 是 $M'$ 的初态),

则令 $A1 \longrightarrow a$ , 有

$L(G') = L(M') = L(G) = L(M) = 0(10)^*$



# 第2章:词法分析 (Lexical Analysis)

---

⌘ 2.1 词法分析程序的功能

⌘ 2.2 词法分析器的设计

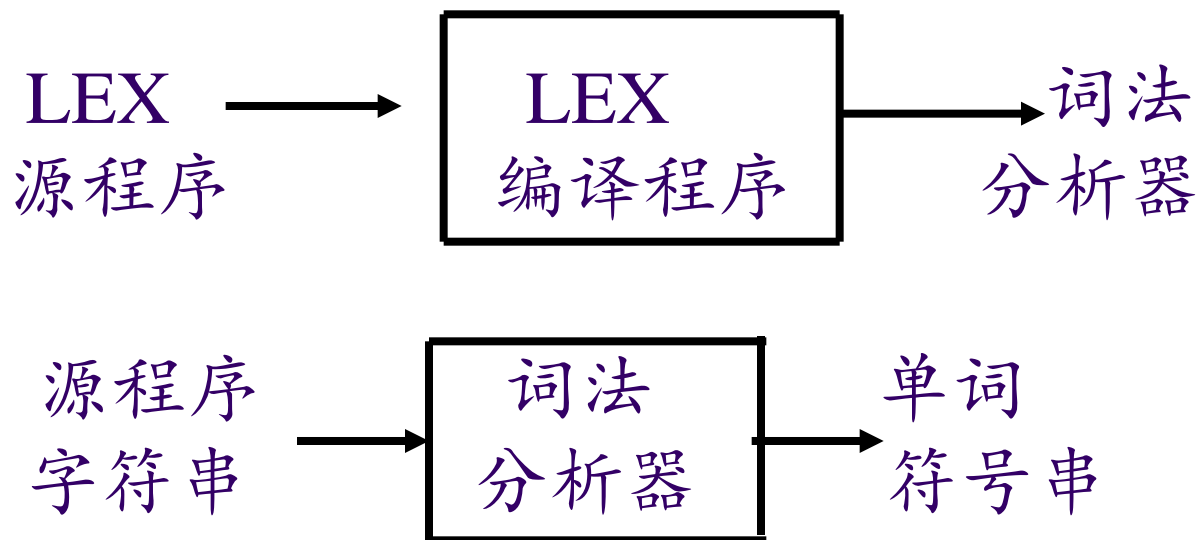
⌘ 2.3 正规表达式 (Regular Expression)

⌘ 2.4 有限自动机

⌘ 2.5 词法分析器的自动生成

## 2.5 词法分析器的自动生成

用LEX语言来写词法分析器，  
LEX编译程序的作用：



## 2.5.1 LEX语言的一般介绍:

LEX程序由两部分生成: (p58~59)

1) 正规式辅助定义式.

2) 词法规则.

一. 辅助定义式:

由一些LEX语句组成, 形式为:

$D1 \longrightarrow R1$

$D2 \longrightarrow R2$

.....

$Dm \longrightarrow Rm$

其中: $R_i$ 为正规式,它定义在

$\Sigma \cup \{D_1, D_2, \dots, D_{i-1}\}$ .

$D_i$  为其简名(用小写字串记)

使用辅助定义式可以方便地定义单词:

如: **FORTRN**中的标识符.

letter  $\longrightarrow$  A|B|C|...|Z

digit  $\longrightarrow$  0|1|2|...|9

iden  $\longrightarrow$  letter|digit

## 二.词法识别规则:

## LEX语句组成:

P1 {A1}

P2                      {A2}

■ ■ ■ ■ ■ ■ ■

$$P_n \quad \{A_n\}$$

其中:  $P_i$ 是词形,由定义在  $\Sigma \cup \{D_1, D_2, \dots, D_m\}$  上的正规式表示.

**A<sub>i</sub>**是动作,当识别出词形**P<sub>i</sub>**后应作的工作。  
词法分析器的功能由**P<sub>i</sub>**和**A<sub>i</sub>**决定.

### 三.词法分析器如何工作:

#### 1) 最长匹配原则:

L扫描输入串,寻找最长的子串匹配某一个 $P_i$ ;并把该子串截下 $\Rightarrow$ TOKEN  
然后调用动作 $A_i$ ,把表示 $P_i$ 的二元式送给语法分析器.

#### 2) 优先匹配原则:

在服从最长匹配的前提下,处于前面的 $P_i$ ,匹配优先权就越高.

### 3) 出错处理:

在输入串中找不到与某一个 $P_i$ 匹配的子串, 要报告出错.

### 4) $A_i$ 返回单词的种别和内部值.在LEX程序中用 RETURN(C, LEXVAL)

标识符: LEXVAL 是TOKEN中的内部值,

常数: LEXVAL是经DTB翻译后在TOKEN中的二进制值。

例题：用LEX识别单词符号： (p58~59)

AUXILIARY DEFINITIONS //正规定义式

Di Ri

.....

RECOGNITION RULES //识别规则

Pi {Ai}

.....



## 2.5.2 超前搜索

例如, Fortran 语句 DO 501 I=1.25

在正规式中引入运算符 “/”,用来表示截断点.  
如,识别基本字DO的规则要写成:

DO/(letter|digit)\*= (letter|digit)\*, {A}

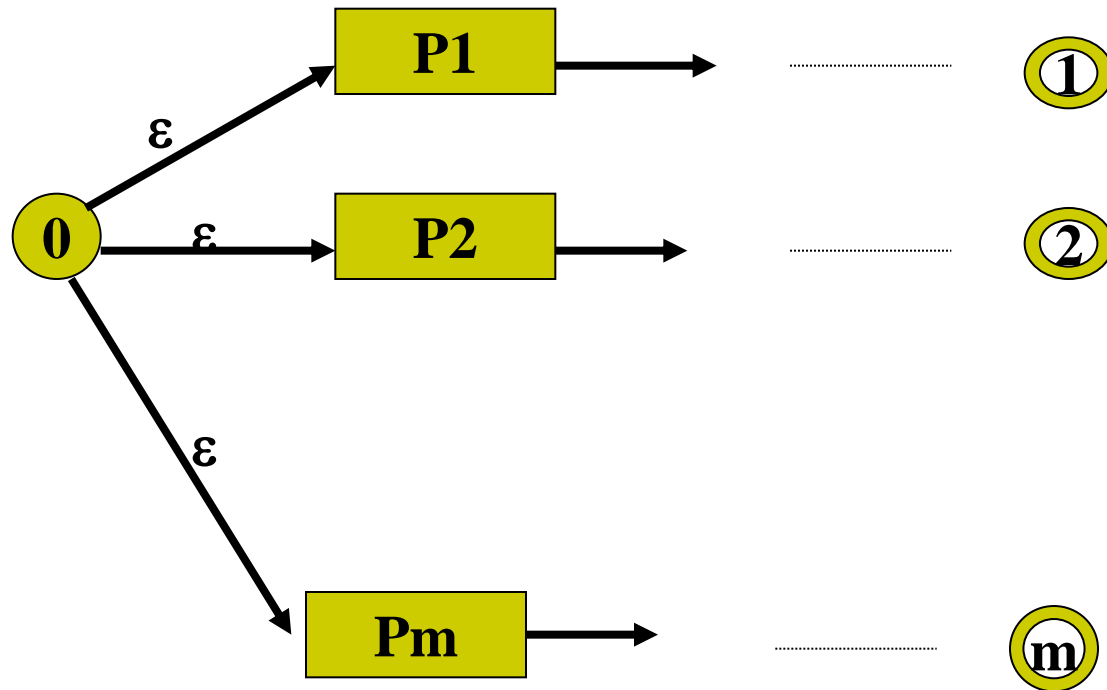
识别时,要求词法分析器L向前扫描到 “,”,从 “/”  
处截断,后一半还给输入串,前一半作为L的输出.

## 2.5.3 LEX的实现

LEX源程序经LEX编译程序翻译成词法分析器L(状态转换表和控制程序组成).

LEX编译过程:

- 1.给每一个 $P_i$ 造相应的NFA  $M_i$ ,
- 2.输入一个初态,经 $\epsilon$ 弧,把所有的NFA  $M_i$ 连接成一个NFA  $M$ .



3.把NFA M 确定化.

4.最小化.

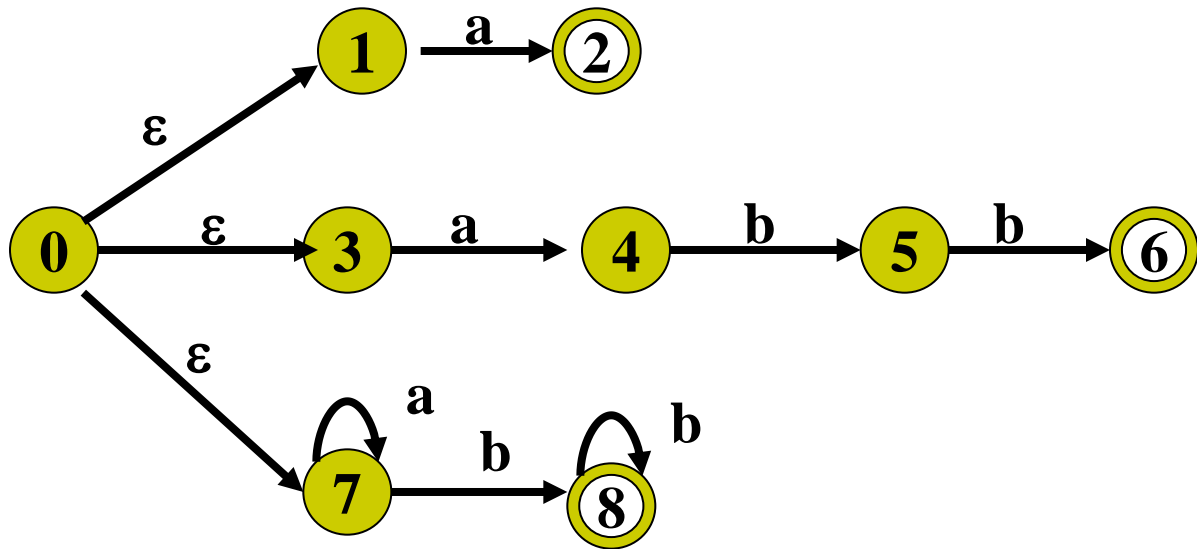
控制程序按状态转换表处理输入串,并能实现最长匹配和最优匹配出错处理.

例如:要识别的词形为:

a            {    }

abb         {    }

a\*bb\*       {    }



确定化后得:

状态	a	b	到达终态时 所识别的词形
初态 0137	247	8	
终态 247	7	58	a
终态 8	--	8	a*bb*
7	7	8	
终态 58	--	68	a*bb*
终态 68	--	8	abb (因为6在8的前面)

# 实现问题：

- 缓冲区预处理，超前搜索，
- 关键字的处理，符号表的实现
- 查找效率，算法的优化实现
- 词法错误处理

# 小结

- 用正规式编写词法，设置单词种别和属性
- 从单词的描述出发，逐步实现词法分析程序
- 词法分析器的自动生成LEX





# 词法分析技术的其他应用：

- 查询语言，信息检索系统
  - 识别由正规式描述的字符串
- 命令语言
  - 识别命令格式
- 报文的词处理
  - 识别报文格式的词处理
- 应用范围：
  - 数据格式可以用三型文法描述。

# 习题

- 1) 考虑 C 语言十六进制常整数,
  - a) 用正规式描述其词法
  - b) 构造对应的 DFA
- 2) 构造下列正规式的状态图, 并以五元组的形式构造对应的 DFA
  - $a ( a b \mid a b^* a )^* b$

# 习题

- 3) 给出下述文法所对应的正规式

- $S \rightarrow 0 A \mid 1 B$

- $A \rightarrow 1 S \mid 1$

- $B \rightarrow 0 S \mid 0$

## 附：Finite Automata

- A *recognizer* for a language is a program that takes a string  $x$ , and answers “yes” if  $x$  is a sentence of that language, and “no” otherwise.
- We call the recognizer of the tokens as a *finite automaton*.
- A finite automaton can be: *deterministic(DFA)* or *non-deterministic (NFA)*

- This means that we may use a deterministic or non-deterministic automaton as a lexical analyzer.
- Both deterministic and non-deterministic finite automaton recognize regular sets.
- Which one?
  - deterministic – faster recognizer, but it may take more space
  - non-deterministic – slower, but it may take less space
  - Deterministic automata are widely used lexical analyzers.

- First, we define regular expressions for tokens; Then we convert them into a DFA to get a lexical analyzer for our tokens.
  - Algorithm1: Regular Expression  $\rightarrow$  NFA  $\rightarrow$  DFA  
(two steps: first to NFA, then to DFA)
  - Algorithm2: Regular Expression  $\rightarrow$  DFA  
(directly convert a regular expression into a DFA)

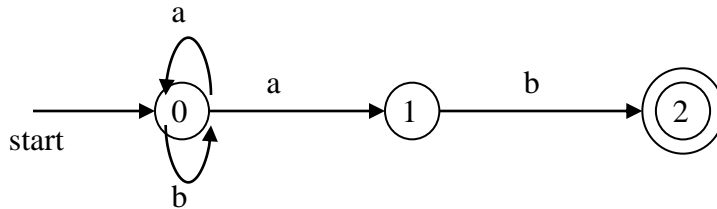
# Non-Deterministic Finite Automaton (NFA)

- A non-deterministic finite automaton (NFA) is a mathematical model that consists of:
  - $S$  - a set of states
  - $\Sigma$  - a set of input symbols (alphabet)
  - move – a transition function move to map state-symbol pairs to sets of states.
  - $s_0$  - a start (initial) state
  - $F$  – a set of accepting states (final states)

- $\varepsilon$ - transitions are allowed in NFAs. In other words, we can move from one state to another one without consuming any symbol.
- A NFA accepts a string  $x$ , if and only if there is a path from the starting state to one of accepting states such that edge labels along this path spell out  $x$ .



# NFA (Example)



Transition graph of the NFA

0 is the start state  $s_0$

$\{2\}$  is the set of final states  $F$

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

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

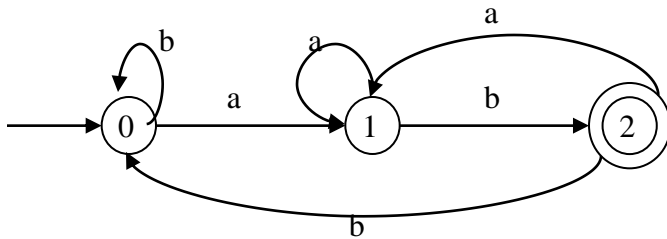
Transition Function:

	a	b
0	$\{0,1\}$	$\{0\}$
1	—	$\{2\}$
2	—	—

The language recognized by this NFA is  $(a|b)^* a b$

# Deterministic Finite Automaton (DFA)

- A Deterministic Finite Automaton (DFA) is a special form of a NFA.
- no state has  $\epsilon$ - transition
- for each symbol  $a$  and state  $s$ , there is at most one labeled edge  $a$  leaving  $s$ .  
i.e. transition function is from pair of state-symbol to state (not set of states)



The language recognized by  
this DFA is also  $(a|b)^* a b$

# Implementing a DFA

- Let us assume that the end of a string is marked with a special symbol (say eos). The algorithm for recognition will be as follows: (an efficient implementation)

```
s ← s0           { start from the initial state }
c ← nextchar      { get the next character from the input string }
while (c != eos) do { do until the end of the string }
  begin
    s ← move(s,c)  { transition function }
    c ← nextchar
  end
if (s in F) then   { if s is an accepting state }
  return “yes”
else
  return “no”
```

# Implementing a NFA

```
S ←  $\epsilon$ -closure( $\{s_0\}$ ) //set all of states can be accessible from  $s_0$  by  $\epsilon$ -transitions
c ← nextchar
while (c != eos) {
    begin
        s ←  $\epsilon$ -closure(move(S,c)) // set of all states can be accessible from a state in S
        c ← nextchar              // by a transition on c
    end
    if ( $S \cap F \neq \Phi$ ) then      { if S contains an accepting state }
        return “yes”
    else
        return “no”
}
```

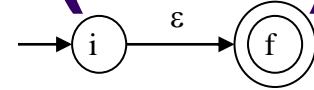
- This algorithm is not efficient.

# Converting A Regular Expression into A NFA (Thomson's Construction)

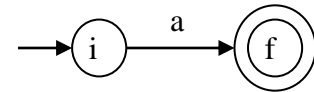
- This is one way to convert a regular expression into a NFA.
- There can be other ways (much efficient) for the conversion.
- Thomson's Construction is simple and systematic method.  
It guarantees that the resulting NFA will have exactly one final state, and one start state.
- Construction starts from simplest parts (alphabet symbols).  
To create a NFA for a complex regular expression, NFAs of its sub-expressions are combined to create its NFA,

# Thomson's Construction (cont.)

- To recognize an empty string  $\varepsilon$

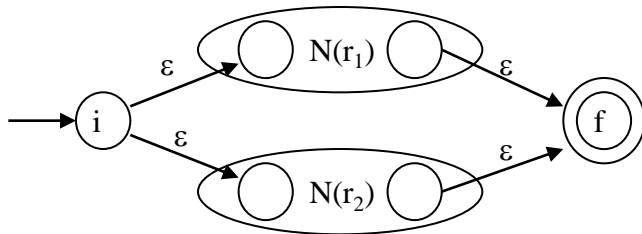


- To recognize a symbol  $a$  in the alphabet  $\Sigma$



- If  $N(r_1)$  and  $N(r_2)$  are NFAs for regular expressions  $r_1$  and  $r_2$

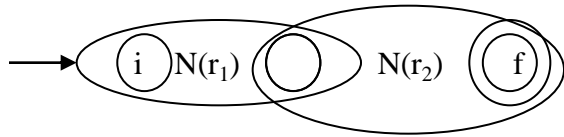
- For regular expression  $r_1 \mid r_2$



NFA for  $r_1 \mid r_2$

# Thomson's Construction (cont.)

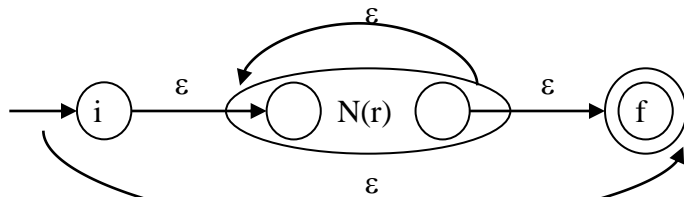
- For regular expression  $r_1 r_2$



Final state of  $N(r_2)$  become final state of  $N(r_1 r_2)$

NFA for  $r_1 r_2$

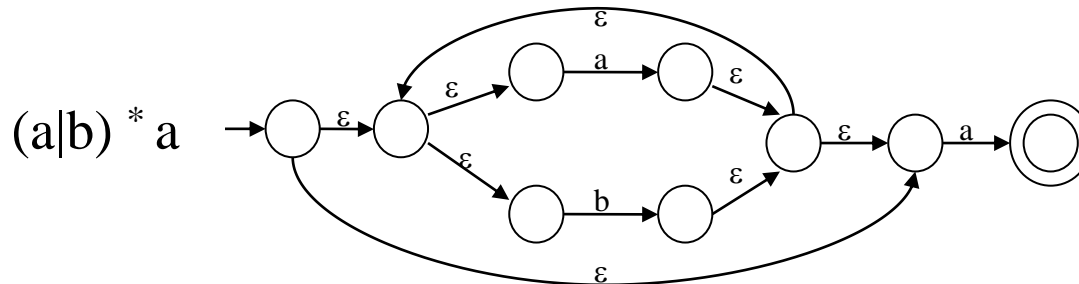
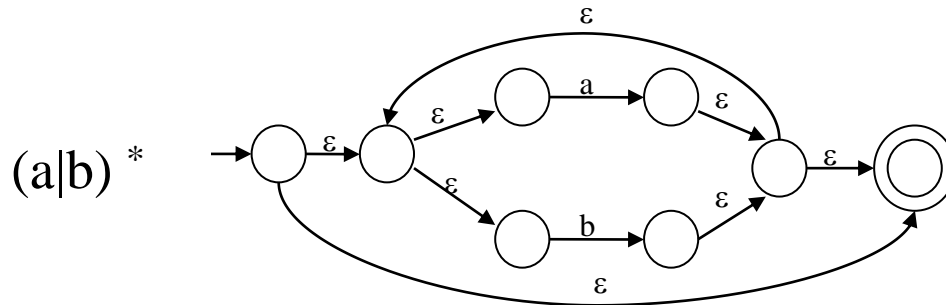
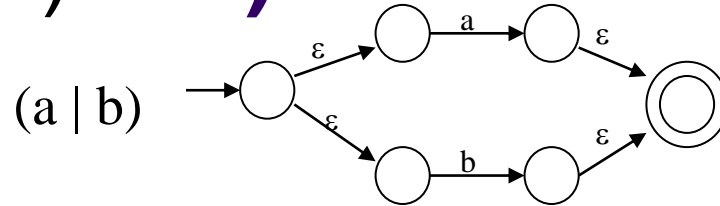
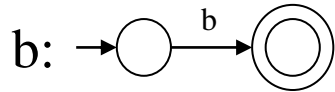
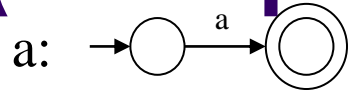
- For regular expression  $r^*$



NFA for  $r^*$

# Thomson's Construction

## (Example - $(a|b)^* a$ )





## Converting a NFA into a DFA (subset construction)

put  $\epsilon$ -closure( $\{s_0\}$ ) as an unmarked state into the set of DFA (DS)

while (there is one unmarked  $S_1$  in DS) do

begin

mark  $S_1$

for each input symbol  $a$  do

begin

$S_2 \leftarrow \epsilon$ -closure(move( $S_1, a$ ))

if ( $S_2$  is not in DS) then

add  $S_2$  into DS as an unmarked state

transfunc[ $S_1, a$ ]  $\leftarrow S_2$

end

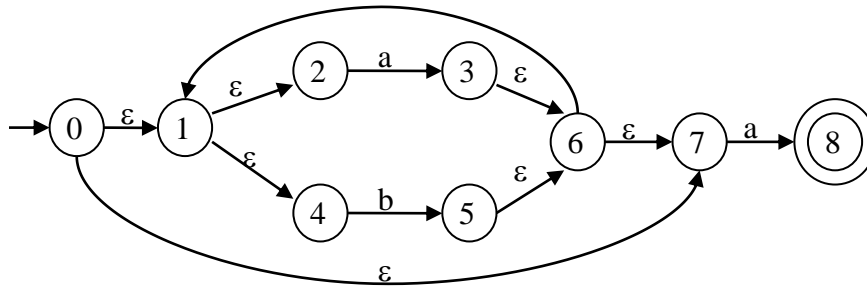
end

$\epsilon$ -closure( $\{s_0\}$ ) is the set of all states can be accessible from  $s_0$  by  $\epsilon$ -transition.

set of states to which there is a transition on  $a$  from a state  $s$  in  $S_1$

- a state  $S$  in DS is an accepting state of DFA if a state in  $S$  is an accepting state of NFA
- the start state of DFA is  $\epsilon$ -closure( $\{s_0\}$ )

# Converting a NFA into a DFA (Example)



$S_0 = \epsilon\text{-closure}(\{0\}) = \{0,1,2,4,7\}$        $S_0$  into DS as an unmarked state

↓ mark  $S_0$

$\epsilon\text{-closure}(\text{move}(S_0, a)) = \epsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = S_1$        $S_1$  into DS

$\epsilon\text{-closure}(\text{move}(S_0, b)) = \epsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = S_2$        $S_2$  into DS

$\text{transfunc}[S_0, a] \leftarrow S_1$        $\text{transfunc}[S_0, b] \leftarrow S_2$

↓ mark  $S_1$

$\epsilon\text{-closure}(\text{move}(S_1, a)) = \epsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = S_1$

$\epsilon\text{-closure}(\text{move}(S_1, b)) = \epsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = S_2$

$\text{transfunc}[S_1, a] \leftarrow S_1$        $\text{transfunc}[S_1, b] \leftarrow S_2$

↓ mark  $S_2$

$\epsilon\text{-closure}(\text{move}(S_2, a)) = \epsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = S_1$

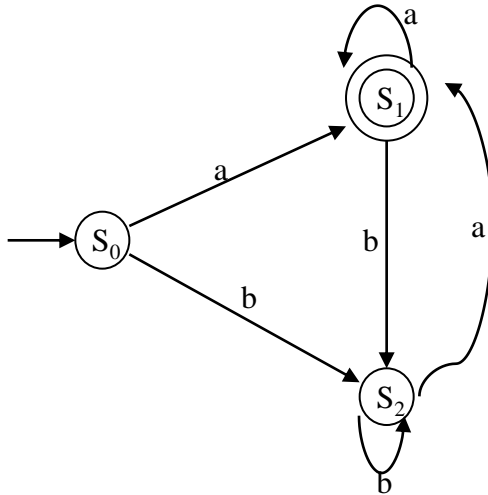
$\epsilon\text{-closure}(\text{move}(S_2, b)) = \epsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = S_2$

$\text{transfunc}[S_2, a] \leftarrow S_1$        $\text{transfunc}[S_2, b] \leftarrow S_2$

# Converting a NFA into a DFA (Example – cont.)

$S_0$  is the start state of DFA since 0 is a member of  $S_0 = \{0, 1, 2, 4, 7\}$

$S_1$  is an accepting state of DFA since 8 is a member of  $S_1 = \{1, 2, 3, 4, 6, 7, 8\}$



# Converting Regular Expressions Directly to DFAs

- We may convert a regular expression into a DFA (without creating a NFA first).
- First we augment the given regular expression by concatenating it with a special symbol #.

$r \rightarrow (r)\#$                       augmented regular expression

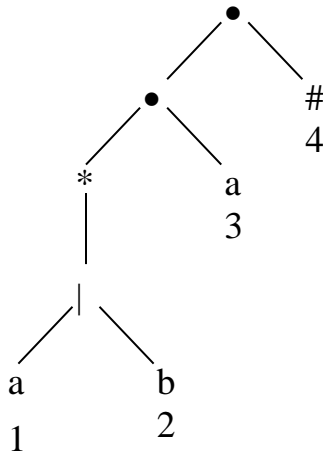
- Then, we create a syntax tree for this augmented regular expression.
- In this syntax tree, all alphabet symbols (plus # and the empty string) in the augmented regular expression will be on the leaves, and all inner nodes will be the operators in that augmented regular expression.
- Then each alphabet symbol (plus #) will be numbered (position numbers).

# Regular Expression $\rightarrow$ DFA

## (cont.)

$(a|b)^* a \rightarrow (a|b)^* a \#$

augmented regular expression



Syntax tree of  $(a|b)^* a \#$

- each symbol is numbered (positions)
- each symbol is at a leaf
- inner nodes are operators