



# 编译原理与设计

北京理工大学 计算机学院

---



# 词法分析：概览

- 两个关键问题
  - 如何定义语言的词法规则
  - 如何识别输入字符串中的单词

Token定义

**正规式**

语言设计和定义

**DFA构造**

Token识别

**确定有限  
状态机**

语言实现



# 词法分析：正规式与正规集

- 设  $\Sigma$  为有限字母表，在  $\Sigma$  上的正规式与正规集可递归定义如下：
  - $\varepsilon$  和  $\Phi$  是  $\Sigma$  上的正规式，它们表示的正规集分别为  $\{\varepsilon\}$  和  $\Phi$ ;
  - 对任何  $a \in \Sigma$ ， $a$  是  $\Sigma$  上的正规式，它的正规集为  $\{a\}$ ;
  - 若  $r, s$  都是正规式，它们的正规集分别为  $R$  和  $S$ ，则  $(r/s)$ 、 $(r \cdot s)$ 、 $(r)^*$  也是正规式，它们分别表示的正规集是： $R \cup S$ ， $RS$ ， $R^*$ 。
- 有限次使用上述三条规则构成的表达式，称为  $\Sigma$  上的正规式，仅由这些正规式表示的集合为正规集。



# 词法分析：正规式与正规集

## ■ 相关说明

- 正规式与相应的正规集是等价的，正规集给出了相应正规式所描述的全部单词（句子）；
  - 正规式的运算结果是正规集；
  - 正规式不是集合，其运算结果正规集是集合， $\Phi$ 是特例；
  - 正规式运算优先级从高到低 “ $( )$ 、 $*$ 、 $\cdot$ 、 $/$ ”；
  - 同级运算从左到右。
-



# 词法分析：正规式与正规集

- 令  $\Sigma = \{0, 1\}$ ，则  $0$ ， $1$ ， $\varepsilon$  和  $\Phi$  是  $\Sigma$  上的正规式；

正规式	正规集
$0 / 1$	$\{0, 1\}$
$0 \cdot 1$	$\{01\}$
$1 \cdot 0$	$\{10\}$
$0^*$	$\{\varepsilon, 0, 00, 000, \dots\}$
$1^*$	$\{\varepsilon, 1, 11, 111, \dots\}$
$(0/1)0^*$	$\{0, 1, 00, 000, \dots, 10, 100, 1000, \dots\}$
$(0/1)01$	$\{001, 101\}$



# 词法分析：正规式与正规集

- 例：令  $\Sigma = \{A, B, 0, 1\}$ 
  - $(A/B)(A/B/0/1)^*$   $\Rightarrow$  { 标识符 }
  - $(0/1)(0/1)^*$   $\Rightarrow$  { 二进制数字串 }
  - $1(01)^* = (10)^*1$

正规式  $r$  所表示的正规集  $R$  是字母表  $\Sigma$  上的语言，称为正则语言，用  $L(r)$  表示，即  $R = L(r)$ 。  $L(r)$  中的元素为字符串，称为  $L(r)$  的句子。

若两个正规式  $r$  和  $s$  所表示的语言  $L(r) = L(s)$ ，则称  $r$ ，  $s$  等价，记为  $r = s$ 。



# 词法分析：正规式与正规集

- C语言有如下单词
    - int、if、else、for、while
    - 标识符、无符号整数
    - <、<=、>、>=、==
  - 对应的正规式描述为
    - *int / if / else / for / while*
    - *<字母>(<字母>|<数字>)\**
    - *</<=/>/>=/=*
-



# 词法分析：正规式与正规集

## ■ 正规式的相关性质

公理/定理	描述
$s / t = t / s$	$/$ 是可交换的
$s / (t / r) = (s / t) / r$	$/$ 是可结合的
$(s t) r = s (t r)$	连接是可结合的
$s (t / r) = s t / s r$	连接对 $/$ 可分配
$(t / r) s = t s / r s$	连接对 $/$ 可分配
$\varepsilon s = s \quad (s \varepsilon = s)$	$\varepsilon$ 是连接的恒等元素
$s^* = (s / \varepsilon)^*$	$*$ 和 $\varepsilon$ 间的关系
$a^* * = a^*$	$*$ 是幂等的





# 词法分析：有限状态自动机





# 词法分析：有限状态自动机

- 确定的有限自动机 (DFA)
    - DFA : **D**eterministic **F**inite **A**utomata
    - 五元组定义:  $M = (S, \Sigma, f, S_0, Z)$ 
      - $S$ : 状态的有限集合, 每个元素  $S_i (S_i \in S)$  称为一个状态
      - $\Sigma$ : 输入字符的有限集合 (或有穷字母表)。每个元素是一个输入字符。
      - $S_0$ :  $M$  的惟一初态 (也称开始状态),  $S_0 \in S$
      - $Z$ :  $M$  的终态集 (或接受状态)  $Z \subseteq S$
      - $f$ : 状态转换函数: 从  $S \times \Sigma \rightarrow S$  的部分映射
-



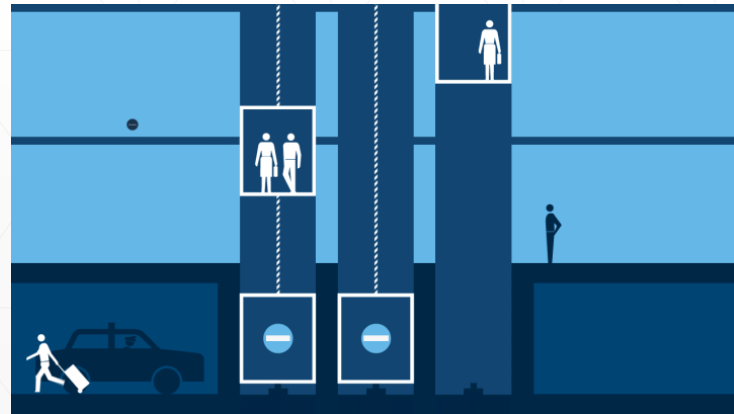
# 词法分析：有限状态自动机

## ■ DFA的说明

- DFA是具有离散输入、输出系统的一个纯数学模型；
- DFA的技巧在于状态的设置；
- DFA映射的唯一性和初态的唯一性。

## ■ 例子

- 计算机系统
- 电梯控制系统





# 词法分析：有限状态自动机

## ■ DFA表示

形式定义

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

状态转换图

$$f(0, a) = 1$$

$$f(0, b) = 2$$

$$f(1, a) = 3$$

$$f(1, b) = 2$$

$$f(2, a) = 1$$

$$f(2, b) = 3$$

状态矩阵

$$f(3, a) = 3$$

$$f(3, b) = 3$$



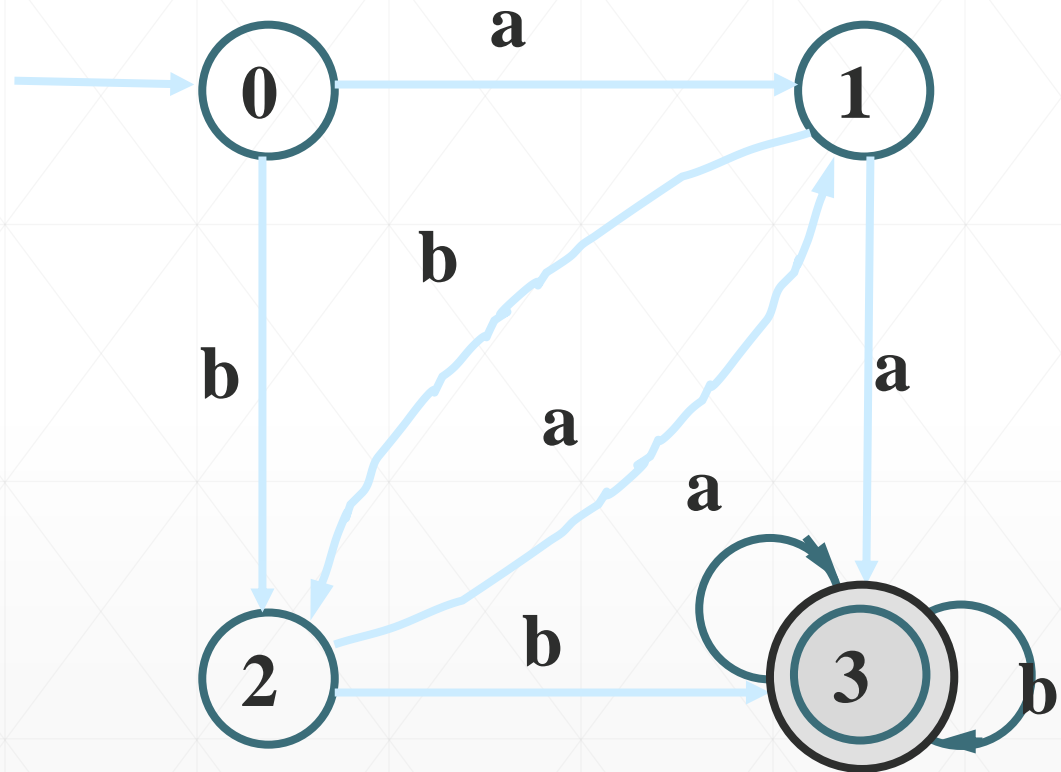
# 词法分析：有限状态自动机

## ■ DFA表示

形式定义

状态转换图

状态矩阵





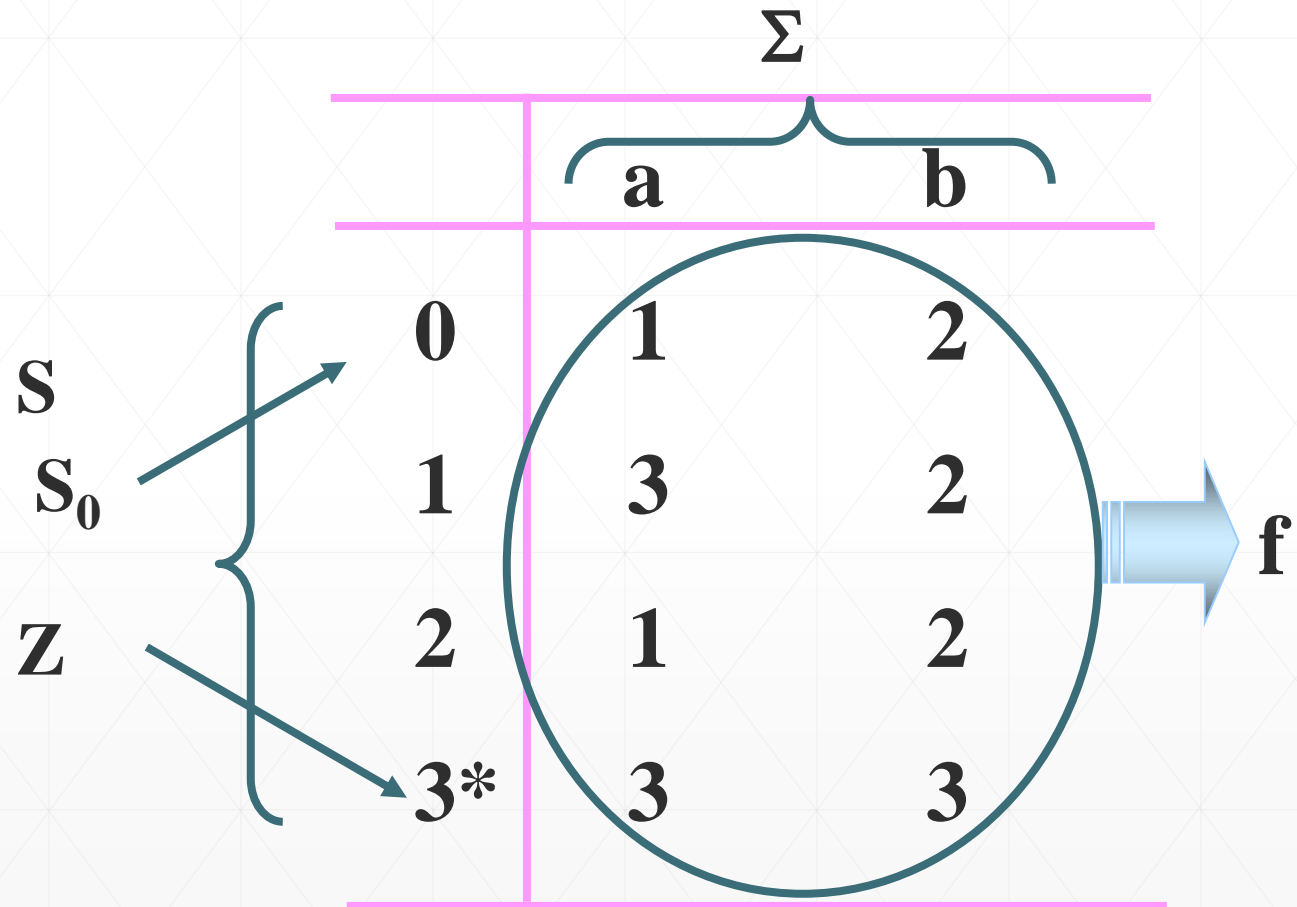
# 词法分析：有限状态自动机

## ■ DFA表示

形式定义

状态转换图

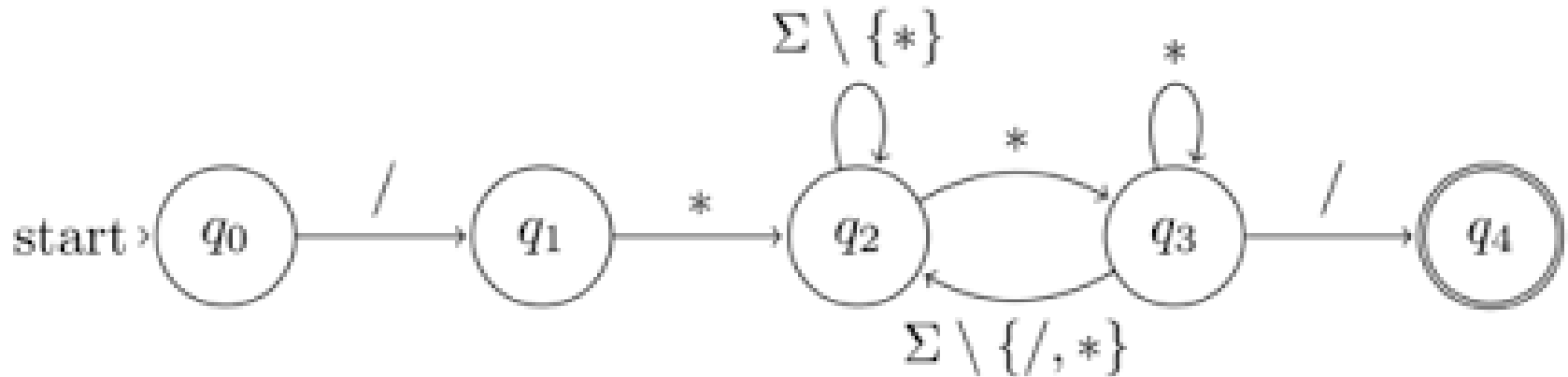
状态矩阵





# 词法分析：有限状态自动机

## ■ DFA: 识别C语言块注释的DFA



```
1 package bit.minisys.minicc.parser;
2
3 import java.util.ArrayList;
4
5 import bit.minisys.minicc.MinicCCfg;
6 import bit.minisys.minicc.util.MinicCUtil;
7
8 /*
9  * PROGRAM --> FUNC_LIST
10  * FUNC_LIST --> FUNC FUNC_LIST | e
11  * FUNC --> TYPE ID '(' ARGUMENTS ')' CODE_BLOCK
12  * TYPE --> INT
13  * ARGS --> e | ARG_LIST
14  * ARG_LIST --> ARG ',' ARG_LIST | ARG
15  * ARG --> TYPE ID
16  * CODE_BLOCK --> '{' STMTS '}'
17  * STMTS --> STMT STMTS | e
18  * STMT --> RETURN_STMT
19  *
20  * RETURN_STMT --> RETURN_EXPR ';'
21  *
22  * EXPR --> TERM EXPR
23  * EXPR --> '+' TERM EXPR | '-' TERM EXPR | e
24  *
25  * TERM --> FACTOR TERM
26  * TERM --> '*' FACTOR TERM | e
27  *
28  * FACTOR --> ID
29  *
30  */
```

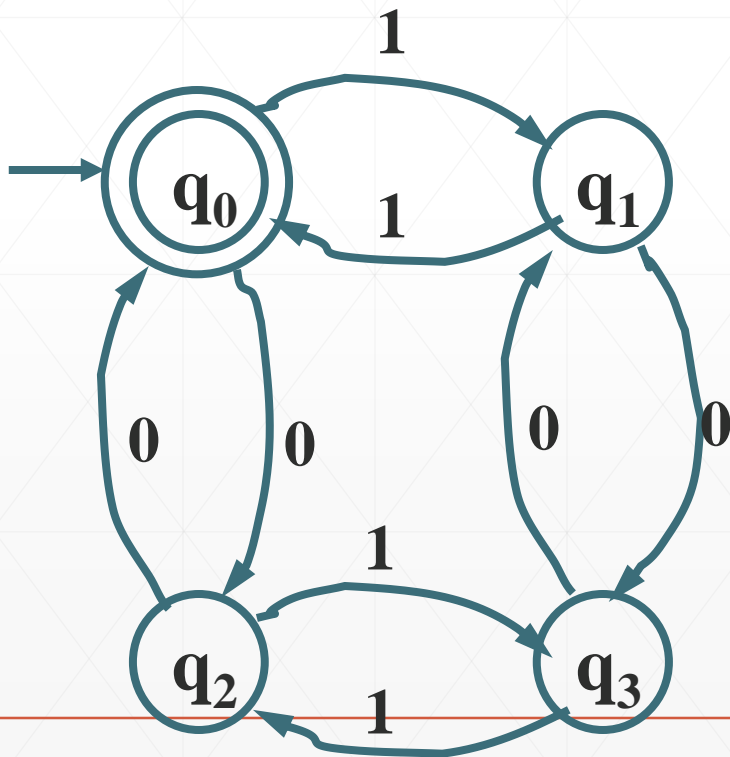
```
//PROGRAM --> FUNC_LIST
public TreeNode program() {
    TreeNode p = new TreeNode(TreeNodeType.TN_TYPE_PROGRAM);
    TreeNode fl = funcList();
    if(fl != null) {
        p.getSubNodes().add(fl);
    }
    return p;
}
```



# 词法分析：有限状态自动机

- DFA: 识别机制

- 例DFA M1可以识别偶数个0或(和)偶数个1。



对1010:

$q_0 \xrightarrow{1} q_1 \xrightarrow{0} q_3 \xrightarrow{1} q_2 \xrightarrow{0} \underline{q_0}$

对11001:

$q_0 \xrightarrow{1} q_1 \xrightarrow{1} q_0 \xrightarrow{0} q_2 \xrightarrow{0} q_0 \xrightarrow{1} \underline{q_1}$





# 词法分析：有限状态自动机

- DFA: 识别机制
    - 对于  $\Sigma$  上的任何字  $\alpha$ ，如果存在一条从初态到某一终态结点的路径，且该路径上所有弧的标记符连接成的字等于  $\alpha$ ，则称  $\alpha$  为 DFA  $M$  所识别(接受)。
    - 若 DFA 仅一个状态结点，该状态结既是初态又是终态，则空字集合  $\{\epsilon\}$  为 DFA  $M$  所接受。
    - 一个 DFA  $M$  所能接受的字的全体记为  $L(M)$ 。
-



# 词法分析：有限状态自动机

- DFA: 语言和等价关系

∴ 有限状态自动机M所接受的语言为:

$$L(M) = \{ \alpha / f(S_0, \alpha) \in Z \ \& \ \alpha \in V_T^* \}$$

设有FA  $M$  和 FA  $M'$  , 若  $L(M) = L(M')$  ,  
则称  $M$  和  $M'$  等价。

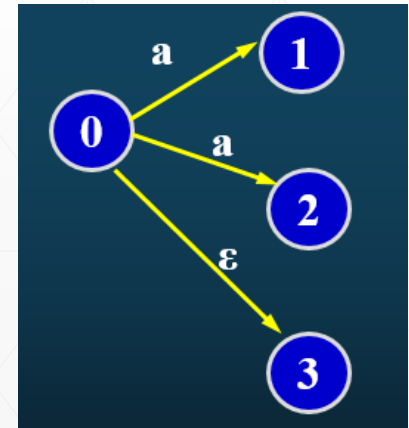


# 词法分析：有限状态自动机

- 实际问题中映射函数往往是多值函数。
  - C语言： `+/++/+=`
- NFA： 定义
  - 非确定的有限自动机M (NFA M) 五元组

$$M = (S, \Sigma, f, S_0, Z)$$

- $S$ ,  $\Sigma$ ,  $Z$ ,  $S_0$ 同DFA
- $F$ ： 状态转换函数：从  $S \times \Sigma^* \rightarrow 2^S$  的映射





# 词法分析：有限状态自动机

## ■ NFA: 表示

形式定义

$$M = (\{q_0, q_1, q_2, q_3, q_4\}, 0, f, q_0, \{q_2, q_4\})$$

$$f(q_0, 0) = q_0 \quad \text{---} \quad f(q_0, 0) = q_3$$

$$f(q_0, 1) = q_0 \quad \text{---} \quad f(q_0, 1) = q_1$$

$$f(q_1, 0) = \Phi \quad \text{---} \quad f(q_1, 1) = q_2$$

$$f(q_2, 0) = q_2 \quad \text{---} \quad f(q_2, 1) = q_2$$

$$f(q_3, 0) = q_4 \quad \text{---} \quad f(q_3, 1) = \Phi$$

$$f(q_4, 0) = q_4 \quad \text{---} \quad f(q_4, 1) = q_4$$

状态转换图

状态矩阵



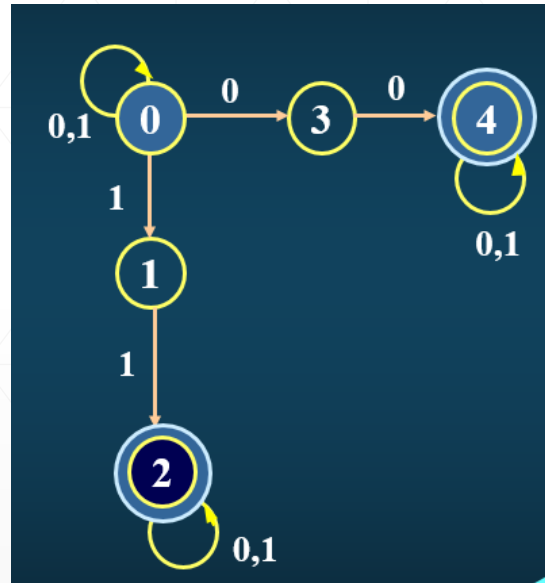
# 词法分析：有限状态自动机

## ■ NFA：表示

形式定义

状态转换图

状态矩阵

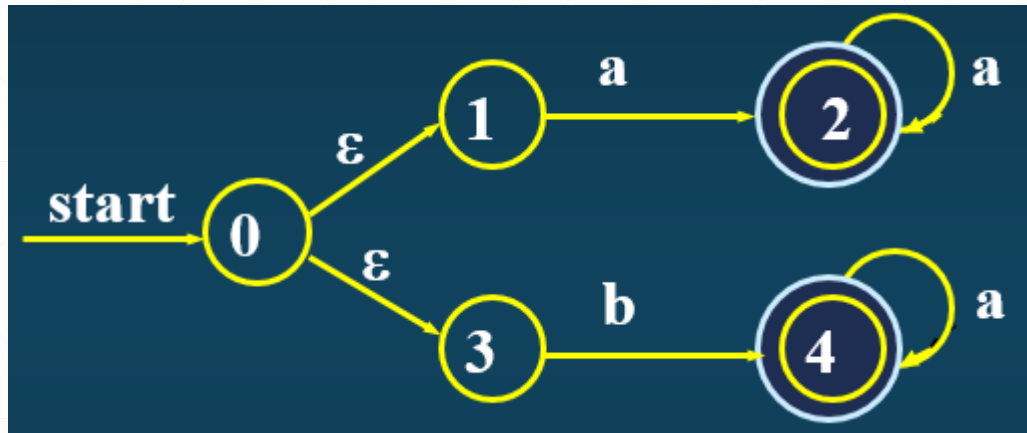


	0	1
0	{0,3}	{0,1}
1		{2}
2*	{2}	{2}
3	{4}	
4*	{4}	{4}



# 词法分析：有限状态自动机

## ■ NFA：识别机制



对字符串aaa的接受路径为0, 1, 2, 2, 2, 接受路径中边的标记是 $\epsilon$ , a, a, a, 它们的连接为字符串aaa,  $\epsilon$ 在连接中消失。



# 词法分析：有限状态自动机

## ■ DFA与NFA区别

DFA	NFA
$\Sigma$	$\Sigma^* (\varepsilon \in \Sigma^*)$
$f(S \times \Sigma)$	$f(S \times \Sigma^*)$
$f(S \times \Sigma) \rightarrow S$	$f(S \times \Sigma^*) \rightarrow 2^S$

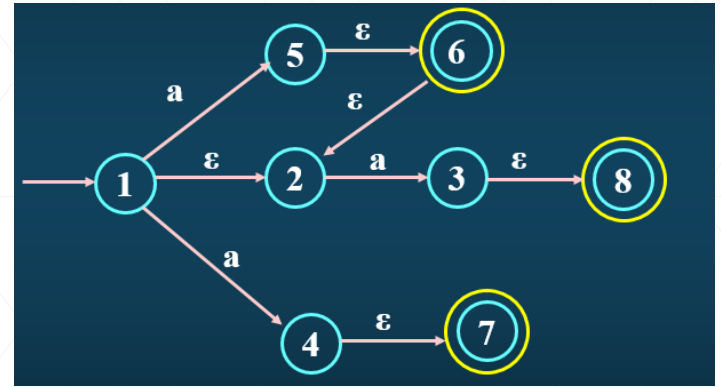
## ■ DFA与NFA等价性

对任何一个NFA  $M$ ，都存在一个DFA  $M'$ ，使  $L(M') = L(M)$ 。



# 词法分析：有限状态自动机

- NFA确定化：子集法
  - 消去 $\epsilon$ 弧： $\epsilon$ -closure(I)
    - $\epsilon$ -closure( $\{5\}$ ) =  $\{5, 6, 2\}$
    - $\epsilon$ -closure( $\{1, 5\}$ ) =  $\{1, 2, 5, 6\}$
  - 解决映射不唯一问题：求 $I_a$ 
    - $I_a = \{2, 5\}_a = \{3, 8\}$
    - $I_a = \{1\}_a = \{2, 3, 4, 5, 6, 7, 8\}$

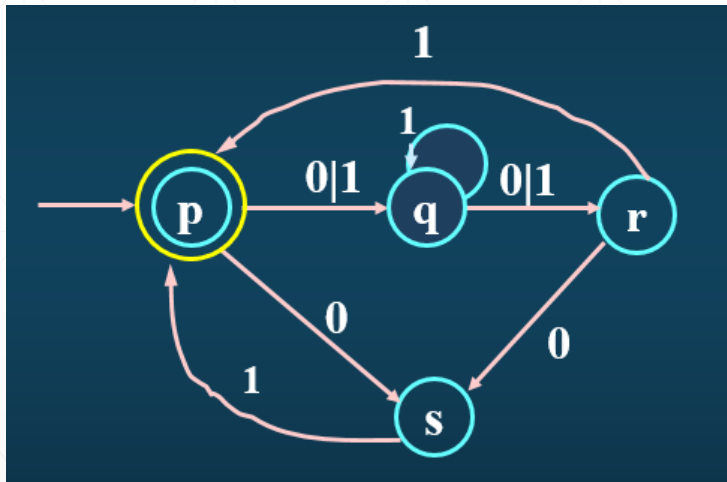






# 词法分析：有限状态自动机

## ■ NFA确定化：子集法



I	I <sub>0</sub>	I <sub>1</sub>
0* {p}	1 {q, s}	{q} 2
1 {q, s}	3 {r}	{q, r, p} 4
2 {q}	3 {r}	{q, r} 5
3 {r}	6 {s}	{p} 0
4* {q, r, p}	7 {q, r, s}	{q, r, p} 4
5 {q, r}	8 {r, s}	{q, r, p} 4
6 {s}		{p} 0
7 {q, r, s}	8 {r, s}	{q, r, p} 4
8 {r, s}	6 {s}	{p} 0



# 词法分析：有限状态自动机

## ■ NFA确定化：子集法

I		I <sub>0</sub>	I <sub>1</sub>
0*	{ p }	1 {q, s}	{ q } 2
1	{q, s}	3 { r }	{ q,r,p } 4
2	{ q }	3 { r }	{ q,r } 5
3	{ r }	6 { s }	{ p } 0
4*	{ q,r,p }	7 { q,r,s }	{ q,r,p } 4
5	{ q,r }	8 { r,s }	{ q,r,p } 4
6	{ s }		{ p } 0
7	{ q,r,s }	8 { r,s }	{ q,r,p } 4
8	{ r,s }	6 { s }	{ p } 0

state	0	1
0*	1	2
1	3	4
2	3	5
3	6	0
4*	7	4
5	8	4
6		0
7	8	4
8	6	0



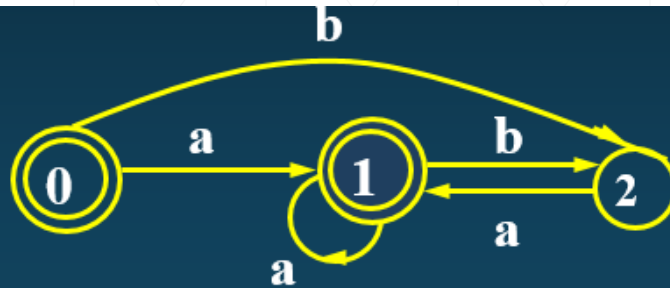
# 词法分析：有限状态自动机

- DFA最小化：划分法
    - **无关状态**：如果从DFA  $M$ 的初态开始，任何输入序列都不能到达的那些状态称为无关状态。
    - **等价状态**：设DFA  $M$ 的两个不同状态  $q_1, q_2$ ，如果对任意输入字符串  $\omega$ ，从 $q_1, q_2$ 状态出发，总是同时到达接收状态或拒绝状态之中，称 $q_1, q_2$ 是等价的。
    - 如果DFA  $M$ 既没有无关状态，且没有彼此等价的状态，则称DFA  $M$ 是**规约的**（即**最小的**DFA  $M$ ）。
-



# 词法分析：有限状态自动机

## ■ DFA最小化：划分法



Step1: 形成初始划分。划分为终态集和非终态集。

$P_0 = (\{0,1\}, \{2\})$

基本划分

考察:  $\{0,1\}_a = \{1\} \subset \{0,1\}$

$\{0,1\}_b = \{2\} \subset \{2\}$

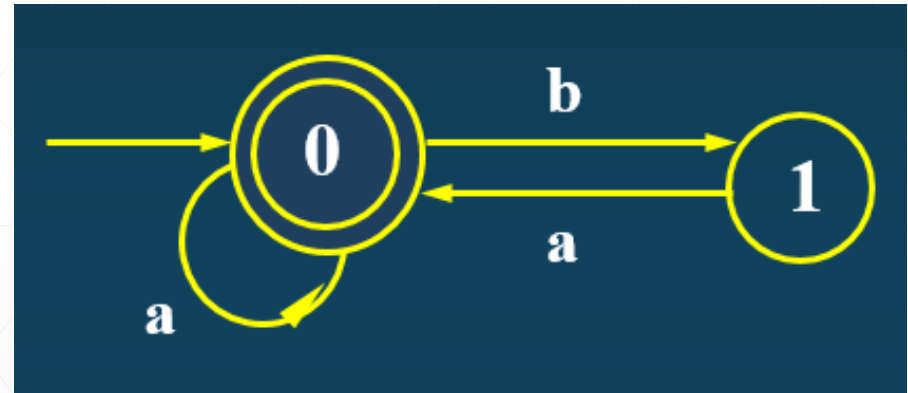
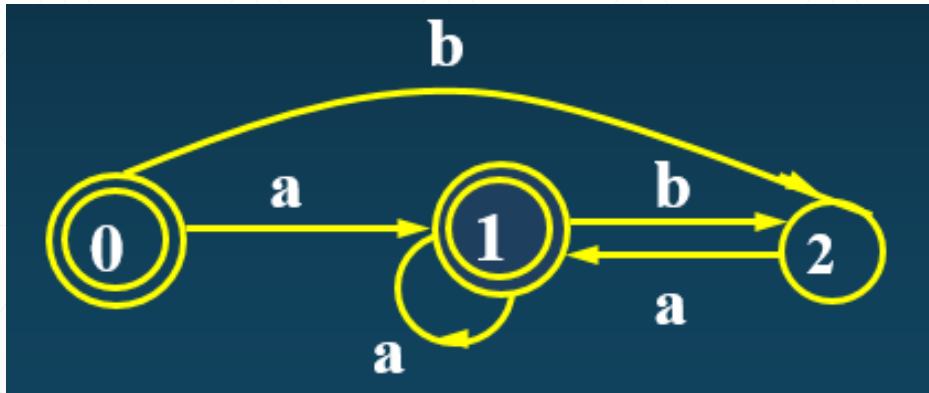
不可对 $\{0,1\}$ 再分

Step2: 重新命名。令 $\{0,1\}$ 为**0**，令 $\{2\}$ 为**1**。



# 词法分析：有限状态自动机

- DFA最小化：划分法





# 词法分析：有限状态自动机

## ■ DFA最小化：划分法

**step1:** 形成初始划分

$\pi = \{ Z, K - Z \}$       //  $K$ 是 $M$ 的所有状态

→ **step2:** 对当前的划分  $\pi = \{ I_1, I_2, \dots, I_m \}$  中的每个状态集  $I_i$  **考察是否可区分**，可区分则进行划分，形成**新的划分**  $\pi_{\text{new}}$ 。

— **step3:** 若  $\pi_{\text{new}} \neq \pi$ ，则将  $\pi_{\text{new}}$  作为  $\pi$  重复**step2**;

**step4:** 对所得的最后划分  $\pi$  重新命名。



# 词法分析：正规式与FA等价性





# 词法分析：正规式与FA等价性

## ■ 定理

- 字母表  $\Sigma$  上的确定的有限自动机  $M$  所接受的语言  $L(M)$  是  $\Sigma$  上的一个正规集;
- 对于  $\Sigma$  上的每一个正规式  $r$ , 存在一个  $\Sigma$  上的非确定有限自动机  $M$ , 使得:  $L(M) = L(r)$ 。

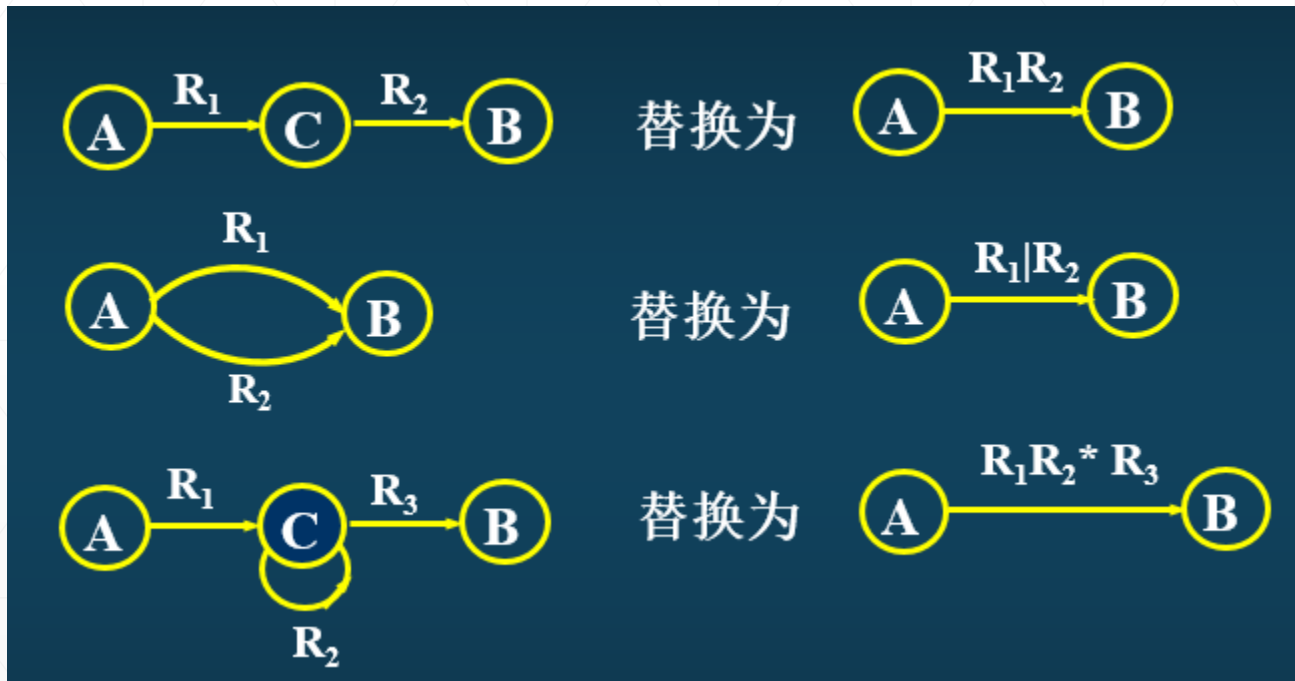
$\Sigma$  上的单词集  $V \in \Sigma^*$  是正规的, 当且仅当存在  $\Sigma$  上的 DFA  $M$ , 使得  $V = L(M)$ 。





# 词法分析：正规式与FA等价性

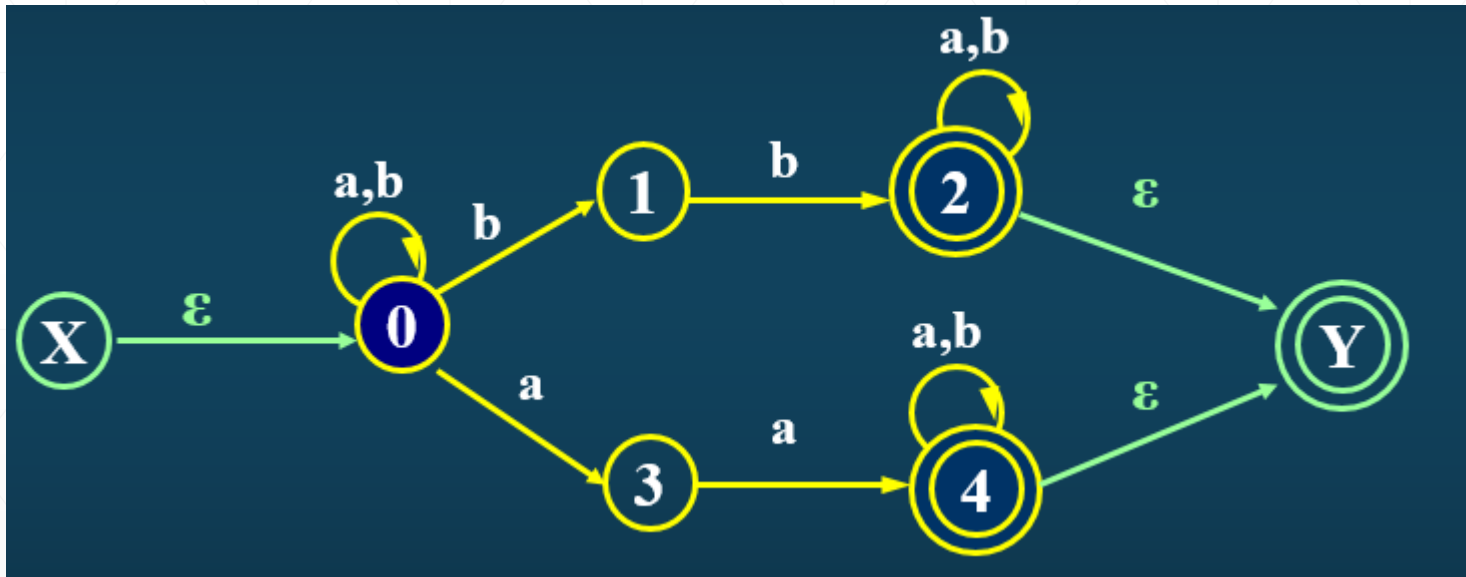
## ■ FA转RE：规则





# 词法分析：正规式与FA等价性

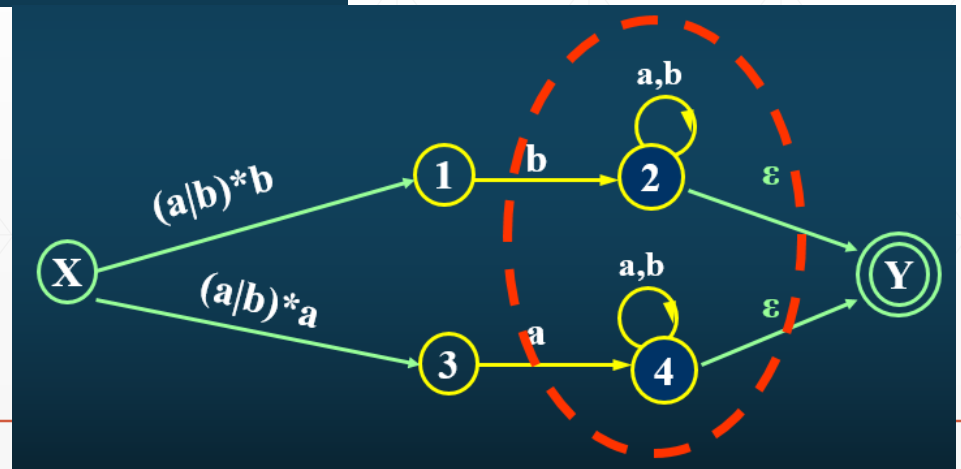
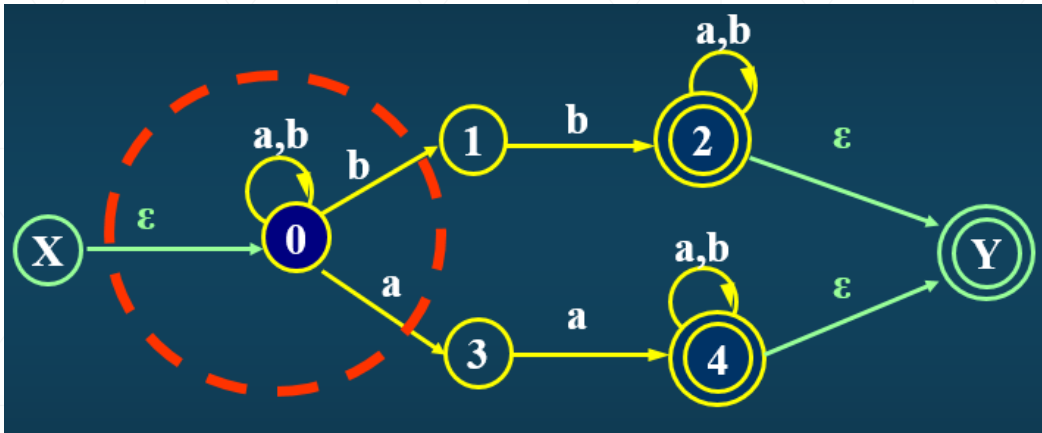
- FA转RE：拓广





# 词法分析：正规式与FA等价性

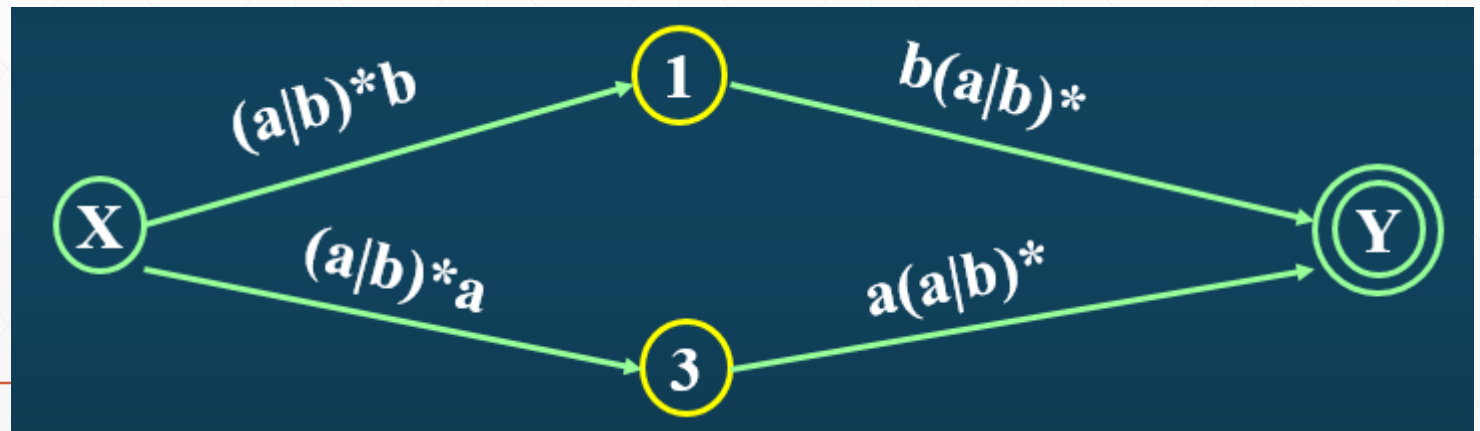
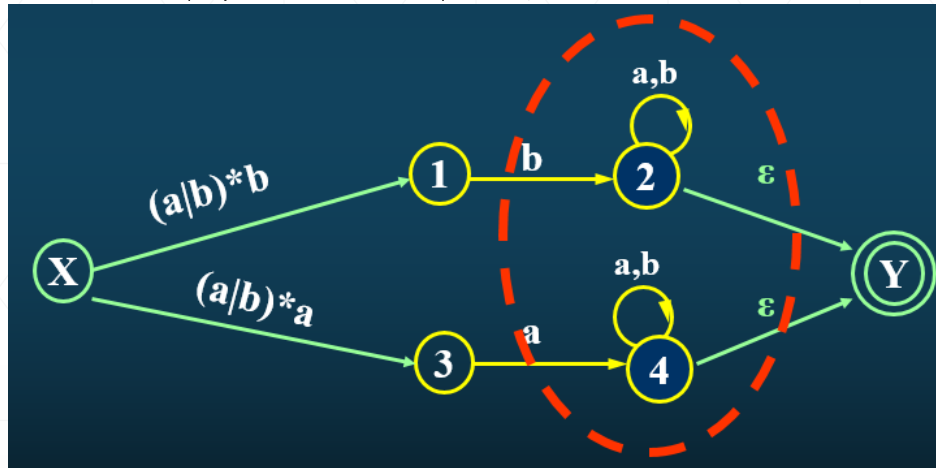
## ■ FA转RE：替换





# 词法分析：正规式与FA等价性

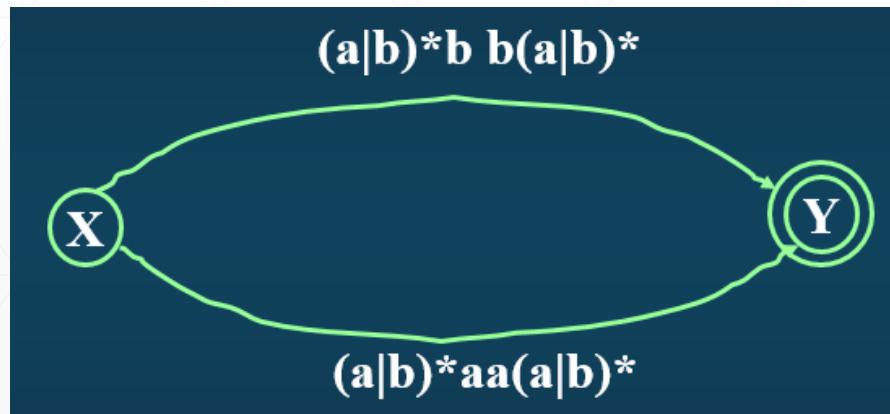
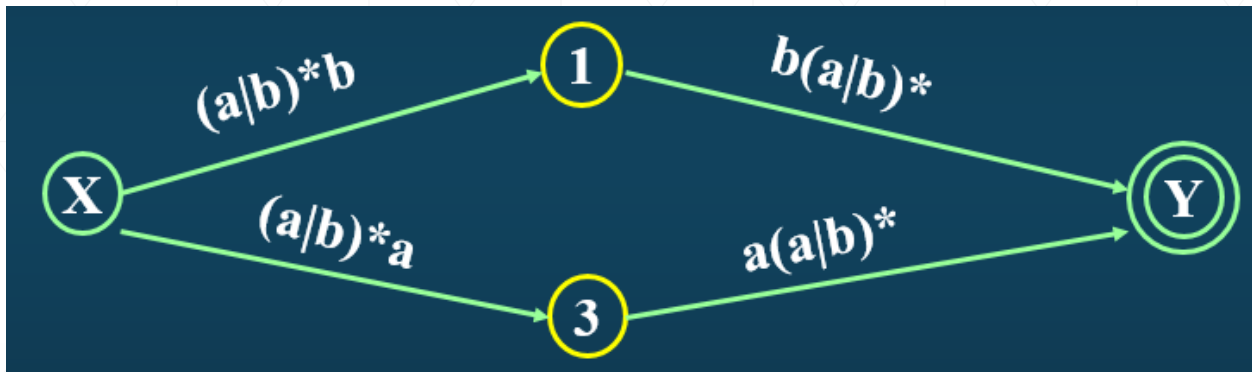
## ■ FA转RE：替换





# 词法分析：正规式与FA等价性

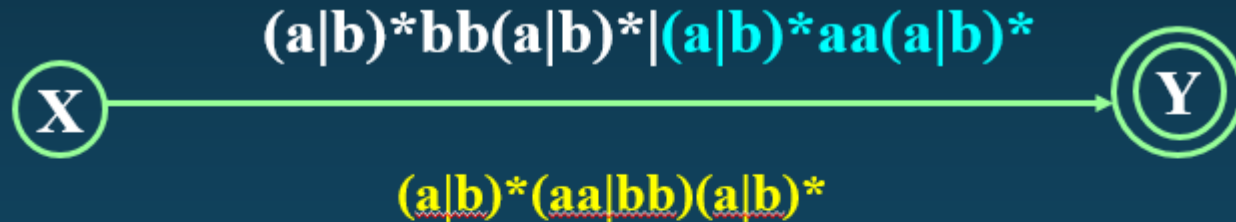
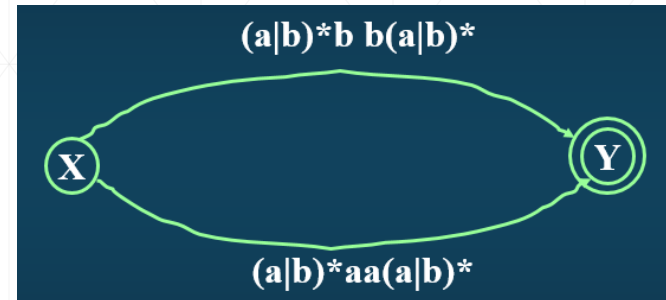
## ■ FA转RE：替换





# 词法分析：正规式与FA等价性

- FA转RE：替换



🔥 注意：(对消弧、消结过程)

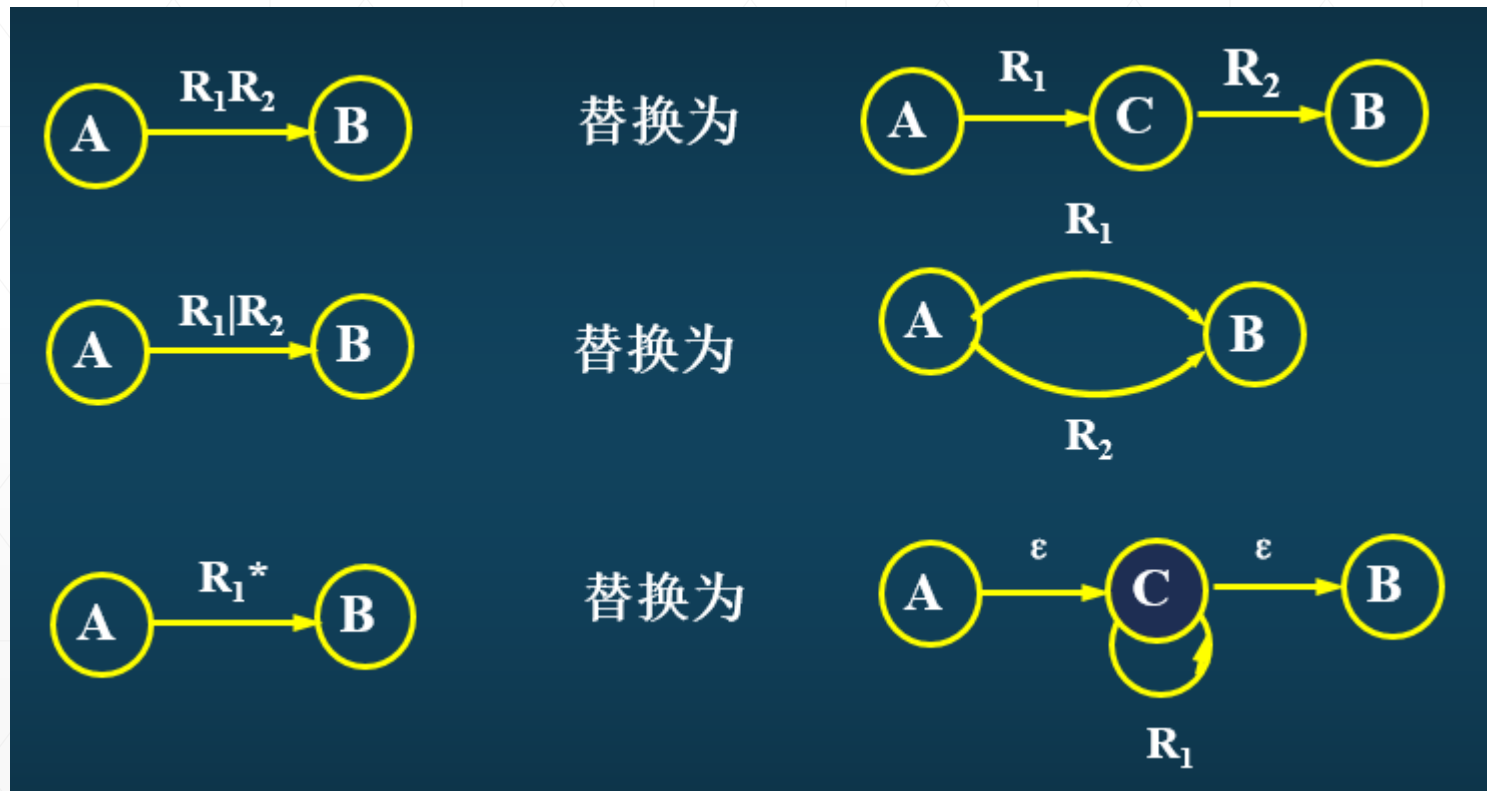
实际是正规式的**合成**过程，其优先级：

$*$  ,  $\cdot$  ,  $|$       ( 高  $\longrightarrow$  低 )



# 词法分析：正规式与FA等价性

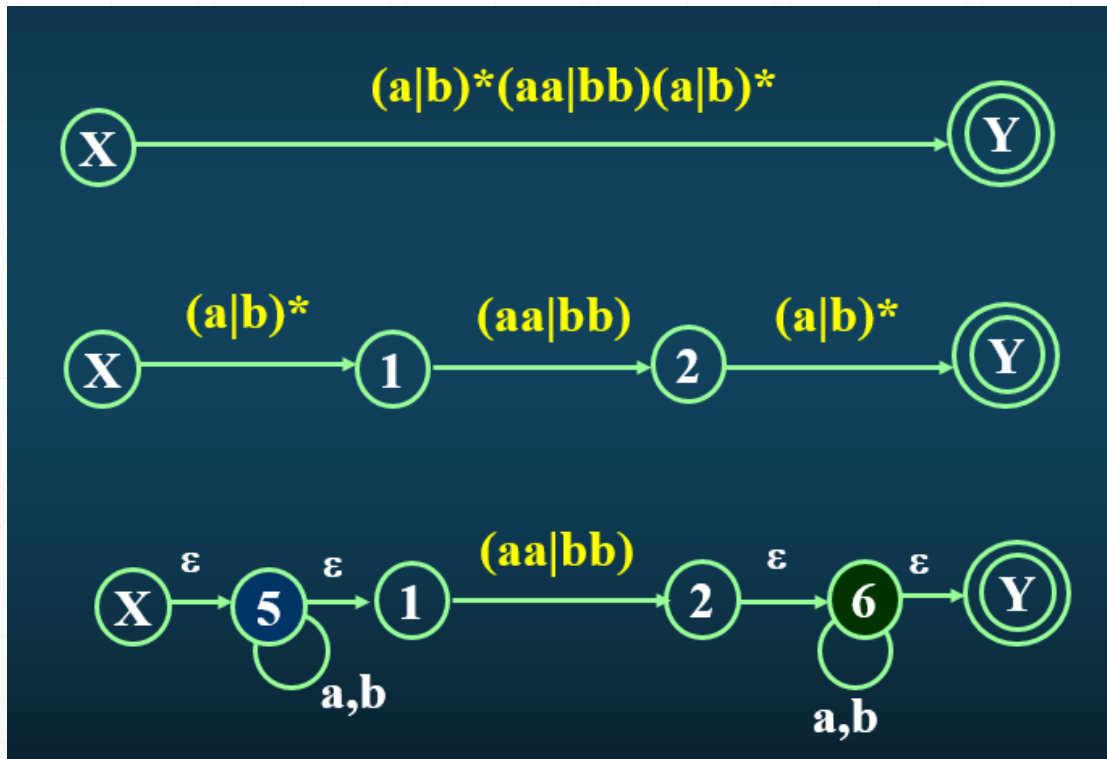
## ■ RE转FA： 规则





# 词法分析：正规式与FA等价性

## ■ RE转FA：拓广与替换

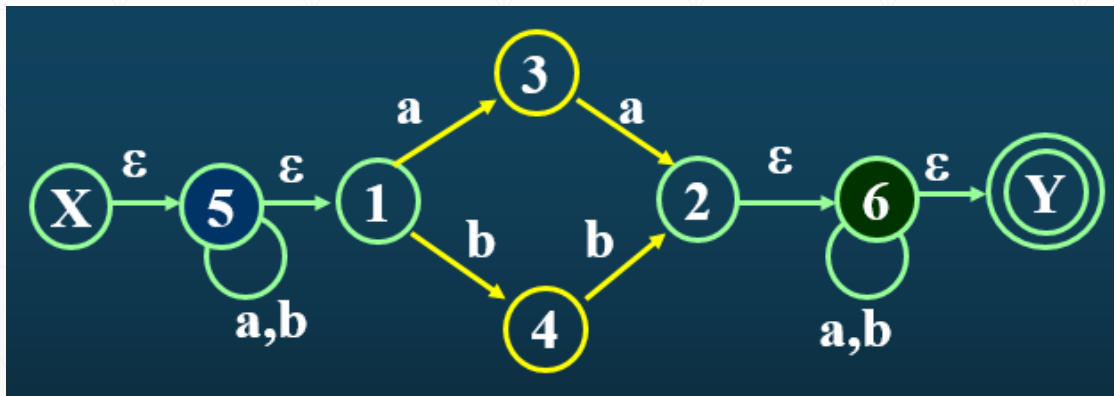
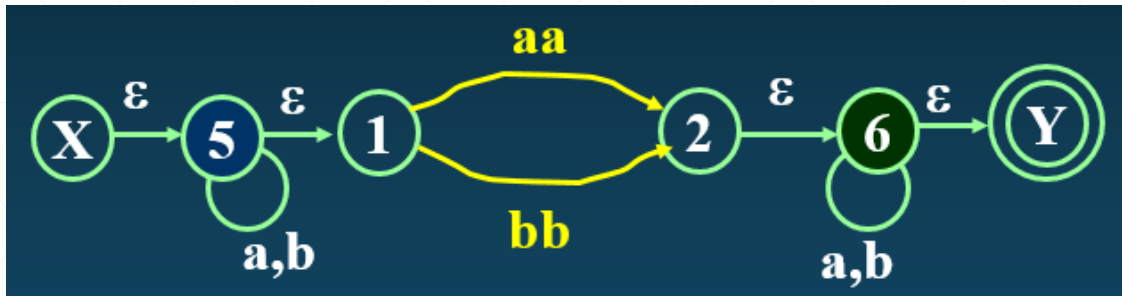






# 词法分析：正规式与FA等价性

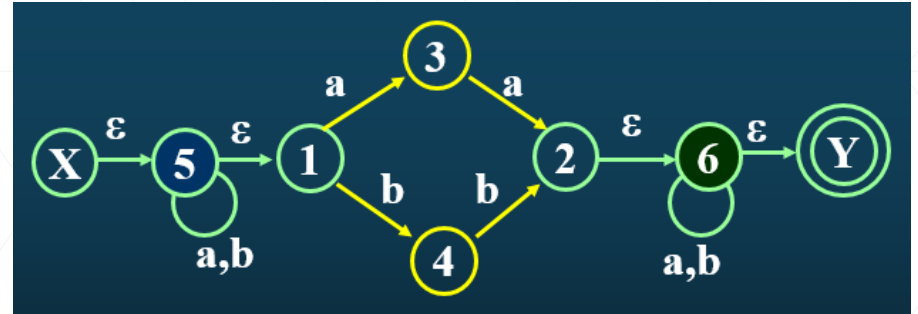
## ■ RE转FA：拓广与替换





# 词法分析：正规式与FA等价性

## ■ RE转FA：确定化



	I	$I_a$		$I_b$	
0	{X,1,5 }	1	{5, 3, 1}	2	{5, 4, 1 }
1	{5, 3, 1 }	3	{5, 3, 1,2,6,Y }	2	{5, 4, 1 }
2	{5, 4, 1 }	1	{5, 3, 1 }	4	{5,4, 1,2,6,Y }
3*	{5, 3, 1,2,6,Y }	3	{5, 3, 1,2,6,Y }	5	{5,4, 1,6,Y }
4*	{5,4, 1,2,6,Y }	6	{5,3, 1,6,Y }	4	{5,4, 1,2,6,Y }
5*	{5,4, 1,6,Y }	6	{5,3, 1,6,Y }	4	{5,4, 1,2,6,Y }
6*	{5,3, 1,6,Y }	3	{5, 3, 1,2,6,Y }	5	{5,4, 1,6,Y }



# 词法分析：正规式与FA等价性

## ■ RE转FA：最小化

state	a	b	state	a	b
0	1	2	0	1	2
1	3	2	1	3	2
2	1	4	2	1	3
3*	3	5	3*	3	3
4*	6	4			
5*	6	4			
6*	3	5			